

```
In [4]: from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
import pandas as pd
import numpy as np
from itertools import cycle, islice
import matplotlib.pyplot as plt
from pandas.plotting import parallel_coordinates
%matplotlib inline
```

```
In [5]: %wd
```

```
Out[5]: 'C:\Users\Wad\l\Desktop\data science analytics\week7 machin learning\Week-7-MachineLearning\weather'
```

```
In [6]: data = pd.read_csv('C:\Users\Wad\l\Desktop\data science analytics\week7 machin learning\Week-7-MachineLearning\weather\weather.csv')
```

```
In [4]: data.shape
```

```
Out[4]: (1587257, 13)
```

```
In [5]: data.head()
```

rowid	hwpren	timestamp	air_pressure	air_temp	avg_wind_direction	avg_wind_speed	max_wind_direction	max_wind_spe
0	0	2011-09-10 00:00:49	912.3	64.76	97.0	1.2	106.0	
1	1	2011-09-10 00:01:49	912.3	63.86	161.0	0.8	215.0	
2	2	2011-09-10 00:02:49	912.3	64.22	77.0	0.7	143.0	
3	3	2011-09-10 00:03:49	912.3	64.40	89.0	1.2	112.0	
4	4	2011-09-10 00:04:49	912.3	64.40	185.0	0.4	260.0	

```
In [6]: data[!data.isnull().any(axis=1)]
```

rowid	hwpren	timestamp	air_pressure	air_temp	avg_wind_direction	avg_wind_speed	max_wind_direction	max_w
0	0	2011-09-10 00:00:49	912.3	64.76	97.0	1.2	106.0	
34790	34790	2011-10-04 10:25:46	915.7	51.08	NaN	NaN	NaN	
35929	35929	2011-10-06 05:24:48	915.2	49.64	NaN	NaN	NaN	
36320	36320	2011-10-06 11:56:49	914.7	50.00	NaN	NaN	NaN	
36321	36321	2011-10-06 11:57:49	914.7	50.00	NaN	NaN	NaN	
36322	36322	2011-10-06 11:58:49	914.7	50.00	NaN	NaN	NaN	
36323	36323	2011-10-06 11:59:49	914.6	50.00	NaN	NaN	NaN	
36324	36324	2011-10-06 12:00:49	914.7	50.00	NaN	NaN	NaN	
36325	36325	2011-10-06 12:01:49	914.6	50.00	NaN	NaN	NaN	
36326	36326	2011-10-06 12:02:49	914.6	50.00	NaN	NaN	NaN	
36327	36327	2011-10-06 12:03:49	914.5	50.18	NaN	NaN	NaN	
36328	36328	2011-10-06 12:04:49	914.5	50.18	NaN	NaN	NaN	
36329	36329	2011-10-06 12:05:49	914.5	50.18	NaN	NaN	NaN	
36330	36330	2011-10-06 12:06:49	914.4	50.18	NaN	NaN	NaN	
36331	36331	2011-10-06 12:07:49	914.4	50.18	NaN	NaN	NaN	
6435	6435	2011-10-25 05:46:50	916.6	51.08	NaN	NaN	NaN	
79098	79098	2011-11-04 04:53:50	911.0	48.92	NaN	NaN	NaN	
79099	79099	2011-11-04 04:55:50	911.0	48.92	NaN	NaN	NaN	
79100	79100	2011-11-04 04:56:50	911.1	48.92	NaN	NaN	NaN	
79101	79101	2011-11-04 04:57:50	911.1	48.92	NaN	NaN	NaN	
79102	79102	2011-11-04 04:58:50	911.1	48.92	NaN	NaN	NaN	
79103	79103	2011-11-04 04:59:50	911.0	48.92	NaN	NaN	NaN	
79104	79104	2011-11-04 05:00:50	911.0	48.92	NaN	NaN	NaN	
79105	79105	2011-11-04 05:01:50	910.9	48.92	NaN	NaN	NaN	
79106	79106	2011-11-04 05:02:50	911.0	48.92	NaN	NaN	NaN	
79107	79107	2011-11-04 05:03:50	910.9	48.92	NaN	NaN	NaN	
79108	79108	2011-11-04 05:04:50	910.9	48.92	NaN	NaN	NaN	
79250	79250	2011-11-04 13:24:50	910.6	48.02	NaN	NaN	NaN	
79609	79609	2011-11-04 19:18:50	908.6	45.14	NaN	NaN	NaN	
79723	79723	2011-11-04 19:19:50	906.9	46.04	NaN	NaN	NaN	
--	--	--	--	--	--	--	--	--
1346164	1346164	2014-03-27 08:46:32	917.1	44.78	NaN	NaN	NaN	
1346165	1346165	2014-03-27 08:47:32	917.1	44.78	NaN	NaN	NaN	
1346166	1346166	2014-03-27 08:48:32	917.1	44.78	NaN	NaN	NaN	
1346167	1346167	2014-03-27 08:49:32	917.1	44.96	NaN	NaN	NaN	
1346168	1346168	2014-03-27 08:50:32	917.1	44.78	NaN	NaN	NaN	
1346169	1346169	2014-03-27 08:51:32	917.2	44.78	NaN	NaN	NaN	
1346170	1346170	2014-03-27 08:52:32	917.1	44.60	NaN	NaN	NaN	
1346171	1346171	2014-03-27 08:53:32	917.1	44.60	NaN	NaN	NaN	
1346172	1346172	2014-03-27 08:54:32	917.1	44.78	NaN	NaN	NaN	
1346173	1346173	2014-03-27 08:55:32	917.1	44.96	NaN	NaN	NaN	
1346174	1346174	2014-03-27 08:56:32	917.1	44.78	NaN	NaN	NaN	
1346175	1346175	2014-03-27 08:57:32	917.0	44.78	NaN	NaN	NaN	
1346176	1346176	2014-03-27 08:58:32	917.0	44.96	NaN	NaN	NaN	
1346177	1346177	2014-03-27 08:59:32	917.1	45.14	NaN	NaN	NaN	
1346178	1346178	2014-03-27 09:00:32	917.1	44.96	NaN	NaN	NaN	
1346179	1346179	2014-03-27 09:01:32	917.2	44.96	NaN	NaN	NaN	
1346180	1346180	2014-03-27 09:02:32	917.1	44.96	NaN	NaN	NaN	
1346181	1346181	2014-03-27 09:03:32	917.2	44.78	NaN	NaN	NaN	
1346182	1346182	2014-03-27 09:04:32	917.2	44.78	NaN	NaN	NaN	
1346183	1346183	2014-03-27 09:05:32	917.2	44.60	NaN	NaN	NaN	
1346184	1346184	2014-03-27 09:06:32	917.3	44.60	NaN	NaN	NaN	
1346185	1346185	2014-03-27 09:07:32	917.3	44.60	NaN	NaN	NaN	
1346186	1346186	2014-03-27 09:08:32	917.3	44.78	NaN	NaN	NaN	
1346187	1346187	2014-03-27 09:09:32	917.4	44.96	NaN	NaN	NaN	
1346188	1346188	2014-03-27 09:10:32	917.3	45.14	NaN	NaN	NaN	
1346189	1346189	2014-03-27 09:11:32	917.4	45.14	NaN	NaN	NaN	
1346190	1346190	2014-03-27 09:12:32	917.4	45.14	NaN	NaN	NaN	
1346191	1346191	2014-03-27 09:13:32	917.4	44.96	NaN	NaN	NaN	
1346192	1346192	2014-03-27 09:14:32	917.5	44.96	NaN	NaN	NaN	
1394844	1394844	2014-04-30 06:51:49	916.7	62.06	NaN	NaN	NaN	

434 rows x 13 columns

```
In [25]: before_rows = data.shape[0]
print(before_rows)
```

1586823

```
In [26]: data = data.dropna()
```

```
In [27]: after_rows = data.shape[0]
print(after_rows)
```

1586823

```
In [28]: before_rows - after_rows
```

```
Out[28]: 0
```

```
In [7]: sampled_df = data[(data['rowID'] % 10) == 0]
```

```
Out[7]: (158726, 13)
```

```
In [8]: sampled_df.head()
```

rowid	hwpren	timestamp	air_pressure	air_temp	avg_wind_direction	avg_wind_speed	max_wind_direction	max_wind_sp
0	0	2011-09-10 00:00:49	912.3	64.76	97.0	1.2	106.0	
10	10	2011-09-10 00:10:49	912.3	62.24	144.0	1.2	107.0	
20	20	2011-09-10 00:20:49	912.2	63.32	63.32	2.0	122.0	
30	30	2011-09-10 00:30:49	912.2	62.60	91.0	2.0	103.0	
40	40	2011-09-10 00:40:49	912.2	64.04	81.0	2.6	88.0	

```
In [9]: sampled_df.describe()
```

	rowid	air_pressure	air_temp	avg_wind_direction	avg_wind_speed	max_wind_direction	max_wind_speed
count	1.58726e+06	158678.000000	158726.000000	458203.937699	0.00	3668612.15	793625.00
mean	7.936209e+05	916.830161	61.851589	162.156100	95.278201	2.775215	163.462144
std	4.582039e+05	0.350171	11.833569	95.278201	95.278201	2.775215	92.452139
min	0.000000e+00	905.000000	51.640000	0.000000	0.000000	0.000000	0.000000
25%	3.968125e+05	914.800000	52.700000	62.000000	1.300000	68.000000	1.600000
50%	7.936209e+05	916.700000	62.240000	182.000000	2.200000	127.000000	2.700000
75%	1.190438e+06	918.700000	70.880000	217.000000	3.800000	223.000000	4.600000
max	1.587260e+06	929.500000	99.500000	359.000000	31.900000	359.000000	36.000000

```
In [10]: sampled_df.describe().transpose()
```

	count	mean	std	min	25%	50%	75%	max
air_pressure	158678.0	793625.000000	458203.937699	0.00	3668612.15	793625.00	1190437.50	1587260.00
air_temp	158726.0	61.851589	11.833569	51.64	95.27	62.24	70.88	99.50
avg_wind_direction	158680.0	162.156100	95.278201	0.00	62.0	182.00	212.00	359.00
avg_wind_speed	158690.0	2.775215	2.057024	0.00	1.3	2.30	3.80	31.90
max_wind_direction	158690.0	163.462144	92.452139	0.00	68.0	187.00	223.00	359.00
max_wind_speed	158690.0	3.400568	2.418027	0.10	1.6	2.70	4.60	36.00
min_wind_direction	158690.0	166.774017	97.441109	0.00	76.0	180.00	212.00	359.00
min_wind_speed	158690.0	2.134664	1.742113	0.00	0.8	1.60	3.00	31.60
rain_accumulation	158725.0	0.000318	0.011236	0.00	0.0	0.00	0.00	3.12
rain_duration	158725.0	0.409627	8.665523	0.00	0.0	0.00	0.00	2960.00
relative_humidity	158725.0	47.609470	26.214409	0.90	24.7	44.70	68.00	93.00

```
In [11]: sampled_df[sampled_df['rain_accumulation'] == 0].shape
```

```
Out[11]: (157812, 13)
```

```
In [12]: sampled_df[sampled_df['rain_duration'] == 0].shape
```

```
Out[12]: (157327, 13)
```

Drop all the Rows with Empty rain_duration and rain_accumulation

```
In [13]: del sampled_df['rain_accumulation']
del sampled_df['rain_duration']
```

```
In [14]: rows_before = sampled_df.shape[0]
sampled_df = sampled_df.dropna()
rows_after = sampled_df.shape[0]
```

```
In [15]: rows_before
```

```
Out[15]: 158726
```

```
In [16]: rows_after
```

```
Out[16]: 157650
```

```
In [17]: rows_before - rows_after
```

```
Out[17]: 46
```

```
In [18]: sampled_df.columns
```

```
Out[18]: Index(['rowid', 'hwpren', 'timestamp', 'air_pressure', 'air_temp', 'avg_wind_direction', 'avg_wind_speed', 'max_wind_direction', 'max_wind_speed', 'min_wind_direction', 'min_wind_speed', 'rain_accumulation', 'rain_duration', 'relative_humidity'], dtype='object')
```

selecxy features of interest for clusterig

```
In [19]: features = ['air_pressure', 'air_temp', 'avg_wind_direction', 'avg_wind_speed', 'max_wind_dir', 'max_wind_speed', 'relative_humidity']
```

```
In [20]: select_df = sampled_df[features]
```

```
In [21]: select_df.columns
```

```
Out[21]: Index(['air_pressure', 'air_temp', 'avg_wind_direction', 'avg_wind_speed', 'max_wind_direction', 'max_wind_speed', 'min_wind_direction', 'min_wind_speed', 'rain_accumulation', 'rain_duration', 'relative_humidity'], dtype='object')
```

```
In [22]: select_df
```

270	911.4	65.92	147.0	0.9	174.0	1.1	36.0
280	911.3	64.76	73.0	1.0	82.0	1.2	43.0
290	911.3	64.84	164.0	1.3	176.0	1.7	43.0
...
158960	914.7	76.46	247.0	0.6	264.0	0.7	43.4
158970	914.6	76.28	208.0	0.7	216.0	0.9	43.7
158980	914.6	76.10	209.0	0.7	216.0	0.9	43.9
158990	914.9	76.28	239.0	0.5	350.0	0.7	43.4
159000	914.9	75.92	344.0	0.4	352.0	0.6	43.9
159010	915.0	75.56	323.0	0.3	348.0	0.5	45.5
159020	915.1	75.56	324.0	1.1	347.0	1.5	46.0
159030	915.1	75.74	1.0	1.3	13.0	1.7	45.8
159040	915.2	75.38	355.0	0.9	1.0	1.1	46.1
159050	915.3	75.38	359.0	1.4	11.0	1.5	46.3
159060	915.4	75.38	11.0	1.1	21.0	1.3	45.7
159070	915.4	75.38	13.0	1.4	24.0	1.6	46.6
159080	915.6	75.20	18.0	1.0	24.0	1.2	46.5
159090	915.6	75.20	356.0	1.7	1.0	1.9	47.2
159100	915.7	75.38	13.0	1.5	24.0	1.7	46.7
159110	915.7	75.02	18.0	1.2	28.0	1.4	46.7
159120	915.7	74.84	25.0	1.4	35.0	1.6	46.5
159130	915.6	74.84	23.0	1.3	30.0	1.5	48.9
159140	915.6	74.84	32.0	1.4	41.0	1.7	45.5
159150	915.6	75.20	23.0	1.1	31.0	1.4	45.7
159160	915.6	75.38	16.0	1.2	28.0	1.5	46.3
159170	915.7	75.38	347.0	1.2	353.0	1.4	48.1
159180	915.6	75.74	326.0	1.2	337.0	1.6	48.3
159190	915.9	75.92	289.0	0.7	309.0	0.9	48.1
159200	915.9	75.74	335.0	0.9	348.0	1.1	47.8
159210	915.9	75.56	330.0	1.0	341.0	1.3	47.8
159220	915.9	75.56	330.0	1.1	341.0	1.4	48.0
159230	915.9	75.56	344.0	1.4	352.0	1.7	48.0
159240	915.9	75.20	359.0	1.3	9.0	1.6	46.3
159250	915.9	74.84	6.0	1.5	20.0	1.9	46.1

158680 rows × 7 columns

```
In [23]: X = StandardScaler().fit_transform(select_df)
Out[23]: array([[ -1.48548281,  0.24544455, -0.68385923, ..., -0.62153592,
    -0.74440309,  0.49233893],
 [ -1.48548281,  0.03247142, -0.19055941, ..., 0.03826701,
    -0.6617126,  0.34710804],
 [ -1.51733107,  0.12374562, -0.65236639, ..., -0.44847286,
    -0.37231683,  0.40639371],
 ...,
 [ -0.39489381,  1.15818654,  1.90856325, ..., 2.63939807,
    -0.78909817,  0.61538018],
 [ -0.36488381,  1.12776181,  2.06599745, ..., -1.67073075,
```