

## Initializing System

```
1 !apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

```
1 !tar xf "/content/drive/My Drive/BigDataAssignment3Files/spark-2.4.5-bin-hadoop2.7.tgz"
2 !pip install -q findspark
```

```
1 import os
2 os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
3 os.environ["SPARK_HOME"] = "/content/spark-2.4.5-bin-hadoop2.7"
```

```
1 import findspark
2 findspark.init()
3 from pyspark.sql import SparkSession
4 from pyspark.context import SparkContext
5 spark = SparkSession.builder.master("local[*]").getOrCreate()
```

## Reading Preprocessed Dataset

```
1 preprocessed_data = spark.read.json("hdfs://udith_gupta_1/processed_data")
```

## Extract Data for EDA

```
1 #Get top 30 crime types in complaints
2 top30_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_crime_list = [ele[0] for ele in top30_crime_type if ele[0] is not None]
```

```
1 #Get top 30 arrests crime types
2 top30_arrests_crime_type = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'A') \
4 .map(lambda row : (row['OFNS_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 arrests_crime_list = [ele[0] for ele in top30_arrests_crime_type if ele[0] is not None]
```

```
1 #Get top 30 location types for crime complaints
2 top30_crime_locations = preprocessed_data.rdd \
3 .filter(lambda row : row['RECORD_TYPE'] == 'C') \
4 .map(lambda row : (row['PREM_TYP_DESC'],1)) \
5 .reduceByKey(lambda key1, key2 : key1 + key2) \
6 .takeOrdered(30,lambda atuple: -atuple[1])
```

```
1 complaints_location_list = [ele[0] for ele in top30_crime_locations if ele[0] is not None]
```

## Generic Imports

```
1 import pandas as pd
2 import plotly.express as px
3 from pyspark.sql.functions import unix_timestamp, from_unixtime
4 from pyspark.sql import functions as F
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import folium
9 from folium.plugins import HeatMap
10
11 %matplotlib inline
```

🔗 /usr/local/lib/python3.6/dist-packages/statsmodels/tools/\_testing.py:19: FutureWarning:  
pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.

## ▼ Time Series Analysis

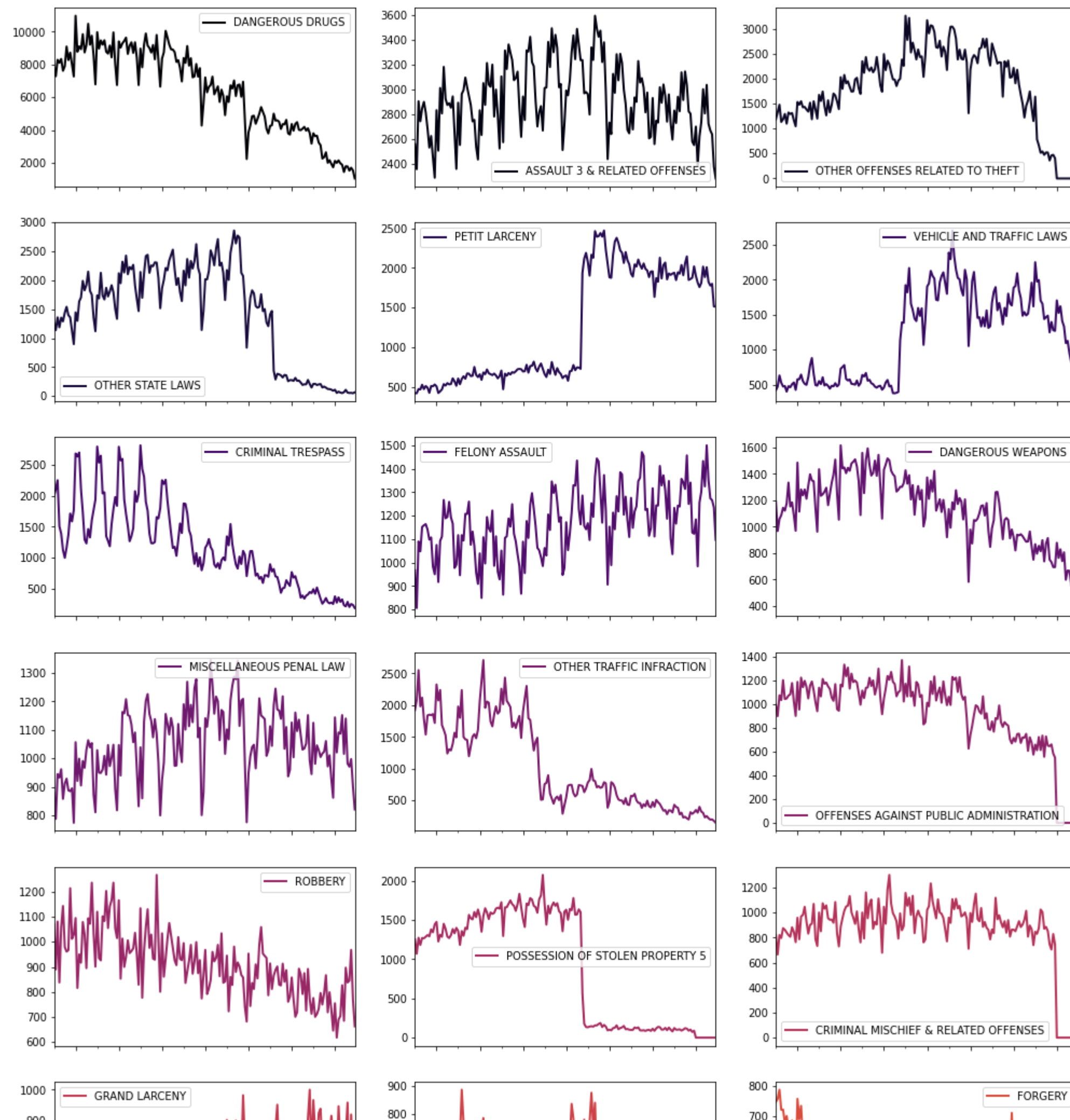
```
1 arrests_df_casted = preprocessed_data.filter(preprocessed_data['RECORD_TYPE'] == 'A').select("*",from_unixtime(unix_timestamp(
```

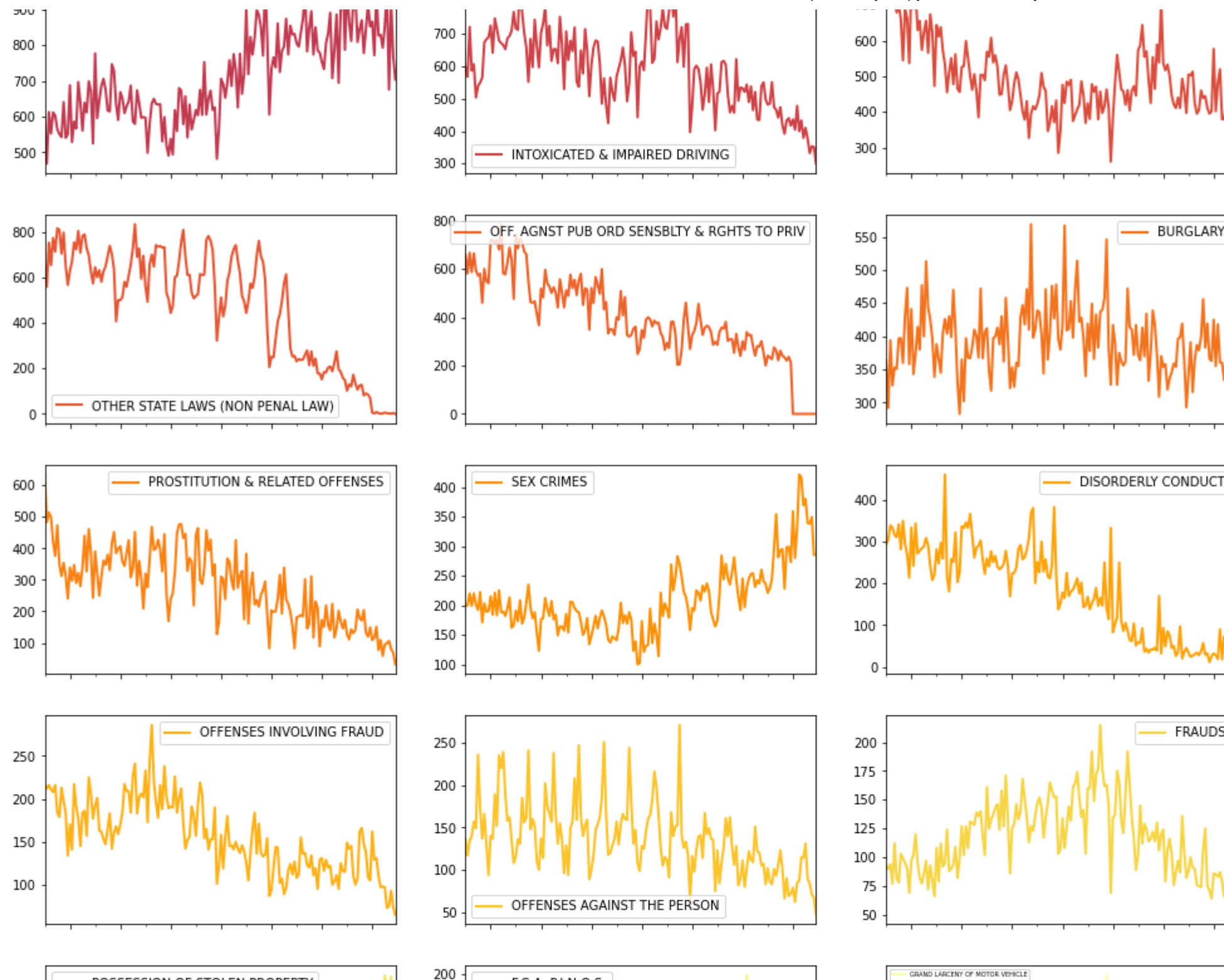
Small Multiples for Major Crimes Types from 2005 to 2018 on monthly grouping

```
1 ` vs Number of Crime for Each Crime Type
2 immarized data to pandas for easy plotting
3 `pe_pivot = arrests_df_casted.groupby('CastedDate').pivot("OFNS_DESC", arrests_crime_list).agg({"*": "count"}).fillna(0).toPanc
4 `pe_pivot['CastedDate'] = pd.to_datetime(day_vs_crimetype_pivot['CastedDate'])
5 `pe_pivot = day_vs_crimetype_pivot.set_index('CastedDate')
```

```
1 #Generate Plot
2 day_vs_crimetype_pivot.resample('M').sum().plot(figsize=(17,40), linewidth=2, cmap='inferno', subplots=True, layout=(-1, 3))
3 plt.legend(prop={'size':5})
4 plt.tick_params(labelsize=8)
```

🔗





```

1 #Helper Functions
2 from sklearn.cluster import AgglomerativeClustering as AC
3
4 def scale_df(df,axis=0):
5     '''
6     A utility function to scale numerical values (z-scale) to have a mean of zero
7     and a unit variance.
8     '''
9     return (df - df.mean(axis=axis)) / df.std(axis=axis)
10
11 def plot_hmap(df, ix=None, cmap='bwr'):
12     '''
13     A function to plot heatmaps that show temporal patterns
14     '''

```

```

15     if ix is None:
16         ix = np.arange(df.shape[0])
17     plt.imshow(df.iloc[ix,:], cmap=cmap)
18     plt.colorbar(fraction=0.03)
19     plt.yticks(np.arange(df.shape[0]), df.index[ix])
20     plt.xticks(np.arange(df.shape[1]))
21     plt.grid(False)
22     plt.show()
23
24 #Plots hmap for scaled data with an additional clustering step for better visuals.
25 def scale_and_plot(df, ix = None):
26     '''
27     A wrapper function to calculate the scaled values within each row of df and plot_hmap
28     '''
29     df_marginal_scaled = scale_df(df.T).T
30     if ix is None:
31         ix = AC(4).fit(df_marginal_scaled).labels_.argsort()
32     cap = np.min([np.max(df_marginal_scaled.as_matrix()), np.abs(np.min(df_marginal_scaled.as_matrix()))])
33     df_marginal_scaled = np.clip(df_marginal_scaled, -1*cap, cap)
34     plot_hmap(df_marginal_scaled, ix=ix)
35     return df_marginal_scaled
36
37 #Normalizes Data with min-max norm
38 def normalize(df):
39     result = df.copy()
40     for feature_name in df.columns:
41         max_value = df[feature_name].max()
42         min_value = df[feature_name].min()
43         result[feature_name] = (df[feature_name] - min_value) / (max_value - min_value)
44     return result
45

```

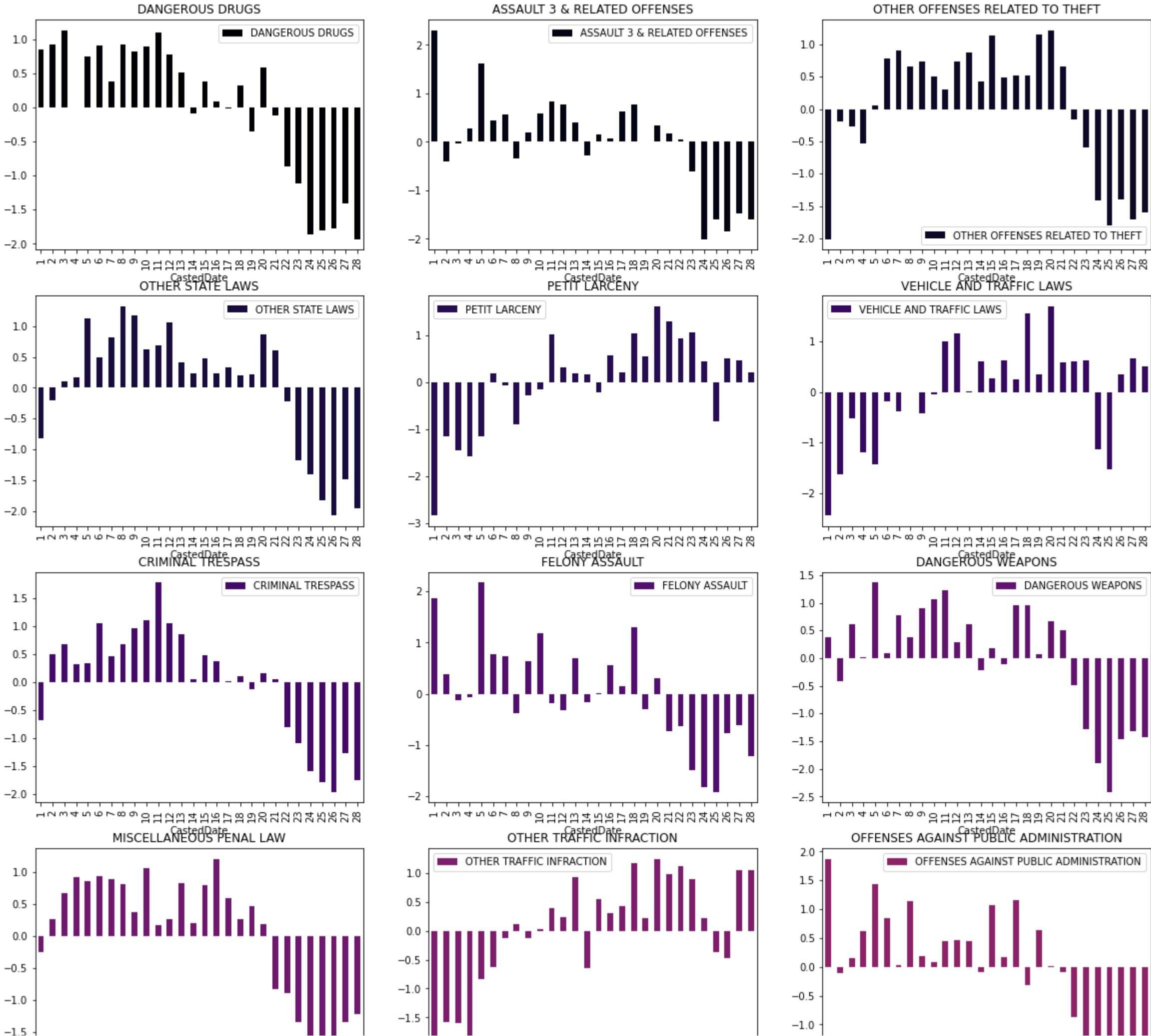
### Small Multiples of Crime Types vs Day of Month

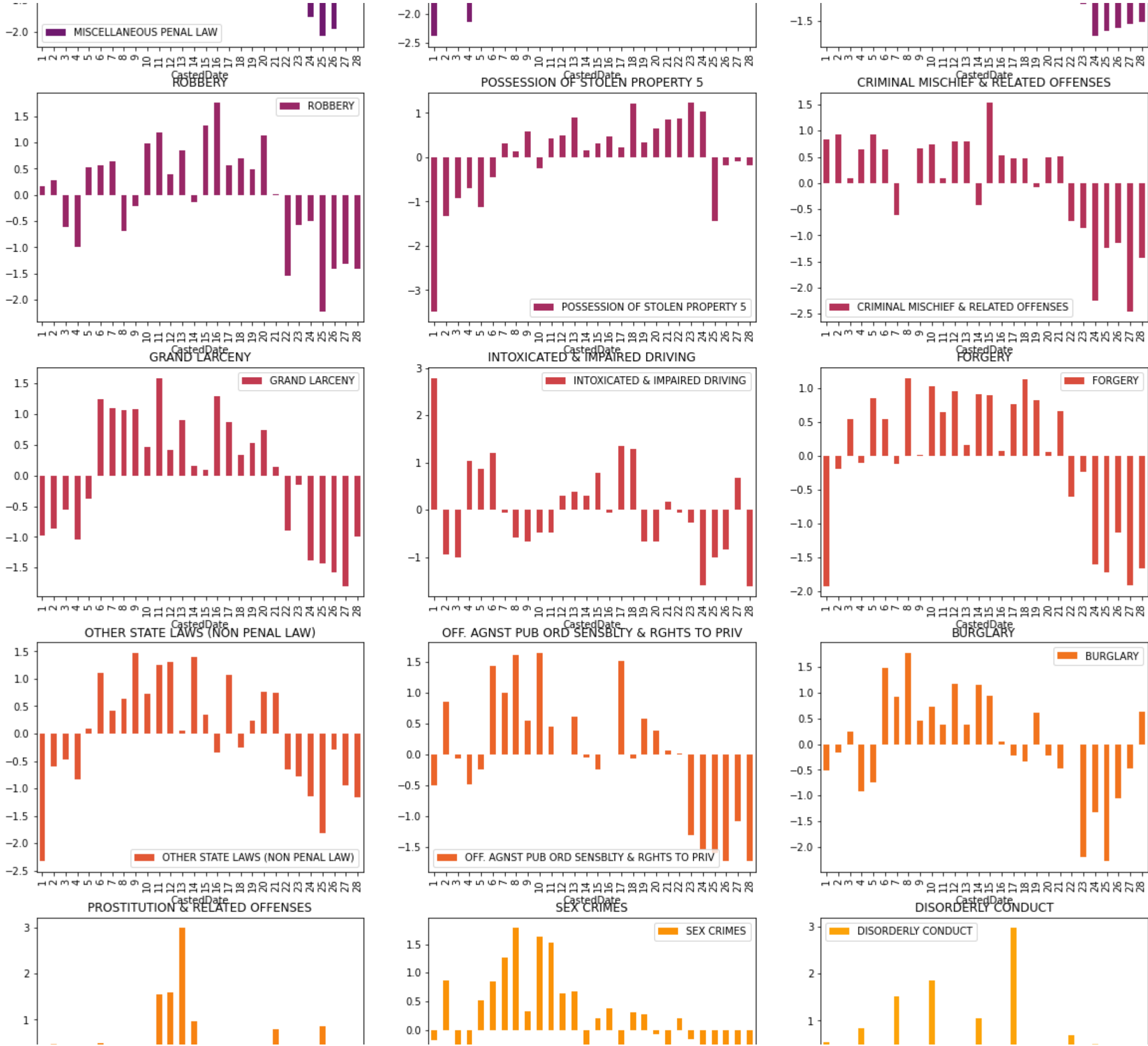
```

1 #Plotting only for first 28 days since 29th, 30th and 31st is not present in Feb
2 scale_df(day_vs_crimetype_pivot.groupby(day_vs_crimetype_pivot.index.day).sum().iloc[:28]).plot(kind='bar', figsize=(20,50), 1
3 plt.legend(prop={'size':5})
4 plt.xlabel('')
5 plt.tick_params(labelsize=5)

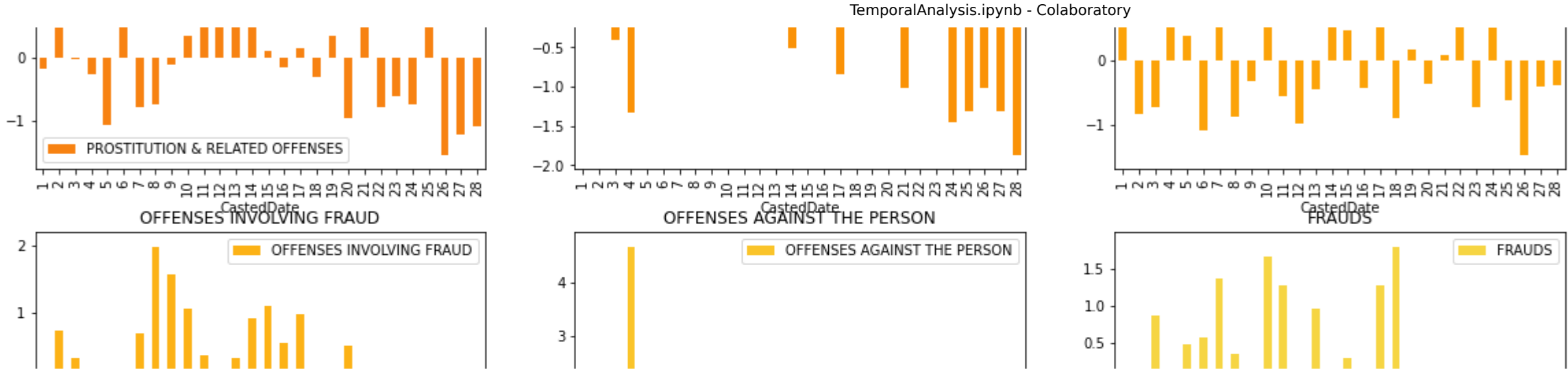
```











Small Multiples of Crime Types vs Month of Year

```
1 cale_df(day_vs_crimetype_pivot.groupby(day_vs_crimetype_pivot.index.month).sum()) \
2 plot(kind='bar', figsize=(22,50), linewidth=2, cmap='inferno', subplots=True, layout=(-1, 3), sharex=False, sharey=False)
3 lt.show()
```





