

# A Machine Learning Strategy for Optimal Stock Trading

Sayan Das                      Alex Nirmal Minj  
dassayan0013@gmail.com      alexnm.math.ug@jadavpuruniversity.in  
Agnik Dasgupta              Debarpan Ghosh  
agnikd76@gmail.com          dghosh4527@gmail.com

October 30, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Aim: . . . . .	3
1.2	Datasets and API used: . . . . .	3
1.3	Github repository and instructions for running the code . . . . .	3
1.4	Packages used: . . . . .	3
<b>2</b>	<b>Methodology</b>	<b>4</b>
2.1	Features and technical indicators to calculate for each stock: . . . . .	4
2.2	Fama-French factors and stock filtering: . . . . .	5
2.3	<i>K</i> -Means clustering and portfolio optimisation: . . . . .	5
2.3.1	Cluster plots for each month . . . . .	6
2.4	Computing portfolio returns: . . . . .	35
<b>3</b>	<b>Final results</b>	<b>36</b>
3.1	Limitations and possible future extensions . . . . .	36
<b>4</b>	<b>References</b>	<b>37</b>

After experimenting with various models and techniques, the framework for this final project was laid down by Sayan, which was then further improvised, refined and fine-tuned by Alex.

Debarpan and Agnik did the tedious work of researching the various datasets and selecting them for training the model, since getting complete usable NIFTY data was next to impossible, we settled on S&P500 data.

We would also like to thank our mentor, Dr. Snehalika Lall, for her invaluable insight.

# 1 Introduction

## 1.1 Aim:

We aim to use an unsupervised machine learning strategy to create a portfolio that outperforms the SPDR S&P 500 (SPY) exchange-traded fund eventually.

## 1.2 Datasets and API used:

- [List of current S&P500 stocks from Wikipedia](#) (as of October 2024)
- [Yahoo Finance API](#)

## 1.3 Github repository and instructions for running the code

The model was coded using Python 3 in a Jupyter notebook. It is available at [this link](#). To run the code given in the `model.ipynb` file, it is of course necessary to have the latest version of Python and JupyterLab Notebook installed. To start the JupyterLab Notebook, open your terminal within the working directory containing the `model.ipynb` file. Within the terminal run the following command:

```
py -m notebook
```

The JupyterLab notebook environment should open up in your browser. Before opening the `model.ipynb` code, go to File> New> Terminal and install the requisite packages:

```
py -m !pip install pandas pandas_ta pandas_datareader numpy matplotlib  
py -m !pip install statsmodels datetime yfinance scikit-learn PyPortfolioOpt
```

## 1.4 Packages used:

- pandas,
- pandas\_ta,
- pandas\_datareader,
- numpy,
- matplotlib,
- statsmodels,
- datetime,
- yfinance,
- scikit-learn,
- PyPortfolioOpt

## 2 Methodology

We now discuss the methodology used.

### 2.1 Features and technical indicators to calculate for each stock:

- Garman-Klass volatility

Defined as

$$r_t = \frac{(\ln H_t - \ln L_t)^2}{2} - 2(\ln 2 - 1)(\ln C_t - \ln O_t)^2,$$

where  $O_t, C_t, H_t, L_t$  denote the opening, closing, high, low prices of day  $t$  respectively.

- Relative strength index (RSI)

Relative strength:

$$RS = \frac{\text{Average gain}}{\text{Average loss}}.$$

Relative strength index:

$$RSI = 100 - \frac{100}{1 + RS}.$$

- Bollinger bands

$$\text{Upper band} = \text{MA}(n) + k\sigma(n),$$

$$\text{Middle band} = \text{MA}(n),$$

$$\text{Lower band} = \text{MA}(n) - k\sigma(n),$$

where  $\text{MA}(n)$  is the moving average over  $n$  periods and  $\sigma(n)$  is the standard deviation over  $n$  periods. Usually we take  $n = 20$  and  $k = 2$ .

- Average true range (ATR)

Given by

$$\frac{1}{n} \sum_{t=1}^n TR(t),$$

where  $n$  is the number of periods, and  $TR(t)$  is the true range of day  $t$  defined by

$$TR(t) = \max(|H_t - L_t|, |H_t - C_{t-1}|, |L_t - C_{t-1}|).$$

- Moving average convergence divergence (MACD)

Given by  $\text{MACD} = \text{EMA}_{12} - \text{EMA}_{26}$ , where  $\text{EMA}_t$  denotes the exponential moving average over  $t$  periods.

- Dollar volume liquidity

Dollar volume liquidity is the price of a stock multiplied by its daily trading volume

## **2.2 Fama-French factors and stock filtering:**

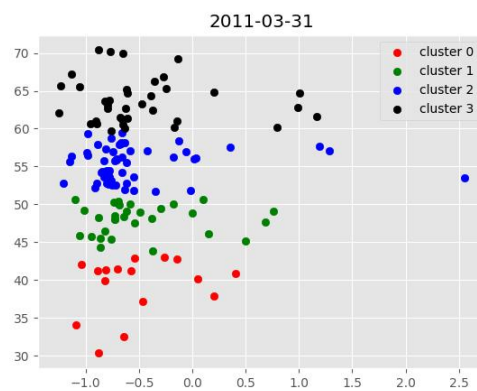
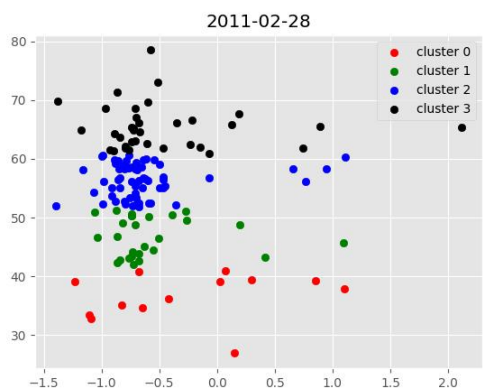
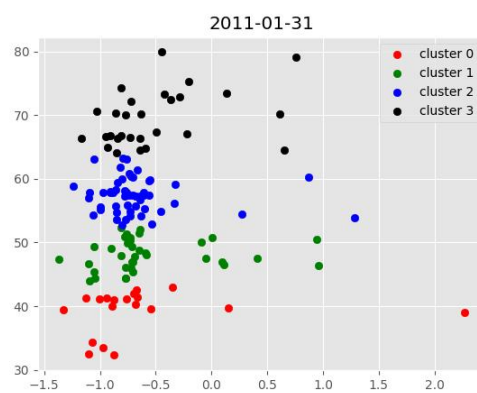
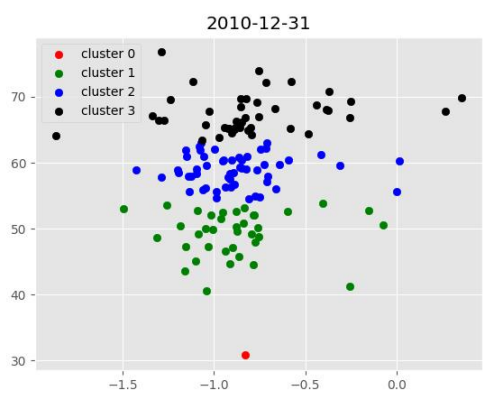
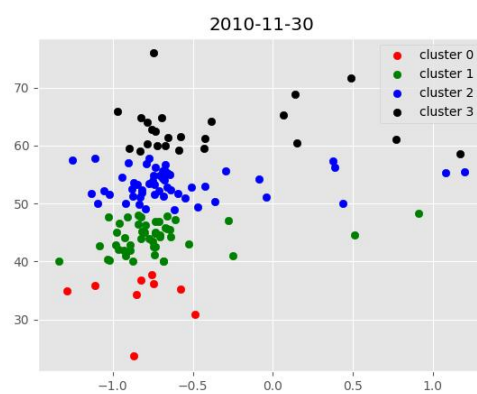
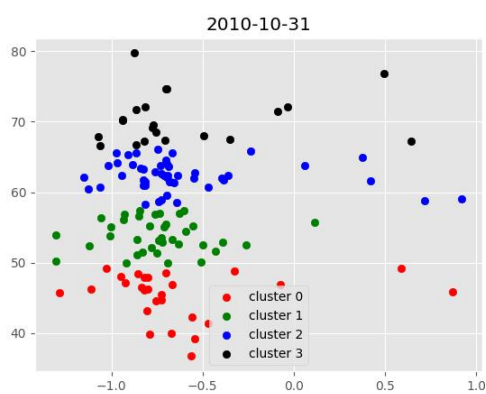
To reduce the model's training time we aggregate the data on a monthly level, filtering the top 150 most liquid stocks per month. We calculate 5-year moving average of dollar volume of each stock before filtering. We then compute monthly returns for different time horizons as features. We download the five Fama-French factors, namely 1. market risk, 2. size, 3. value, 4. operating profitability, and 5. investment

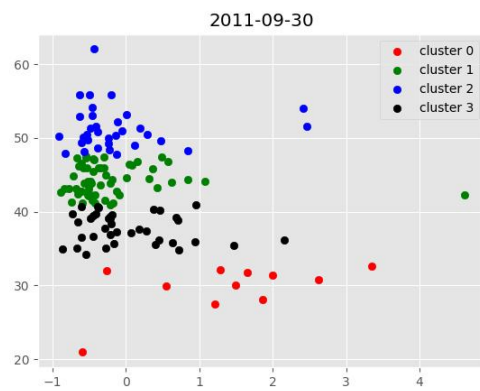
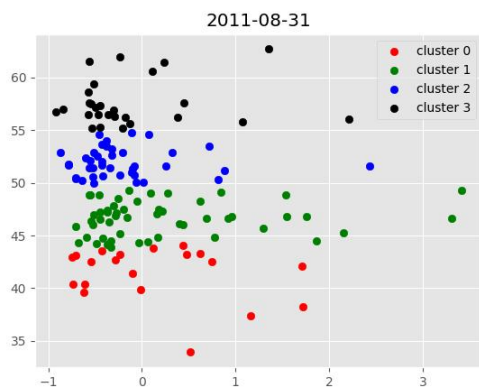
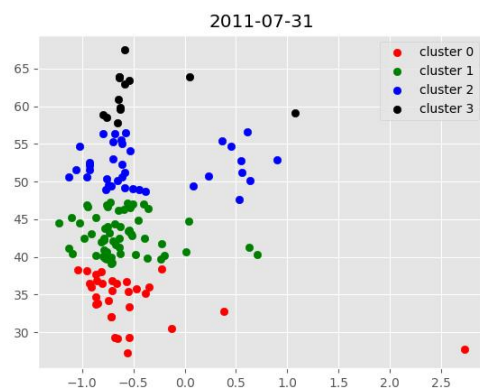
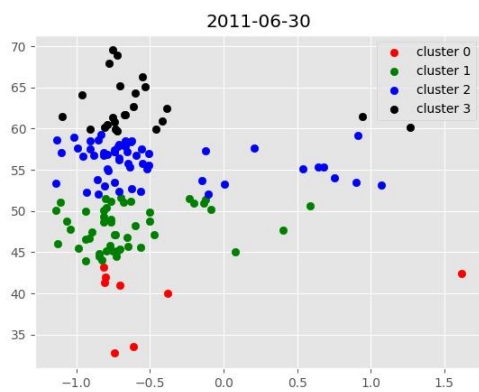
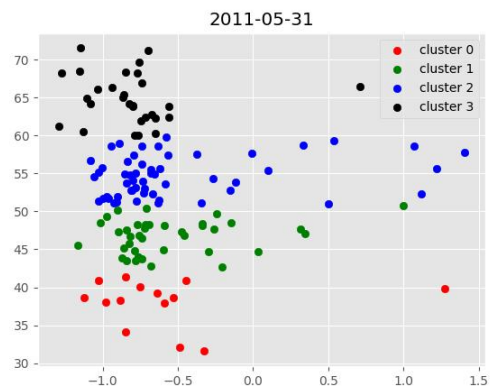
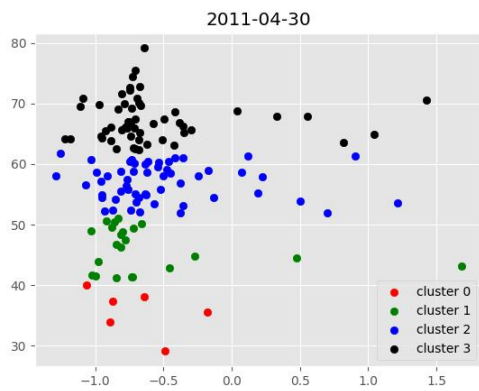
which have been empirically proven for asset returns and are used to assess the risk/return profile of portfolios. It is natural to include past factor exposures as financial features in our model.

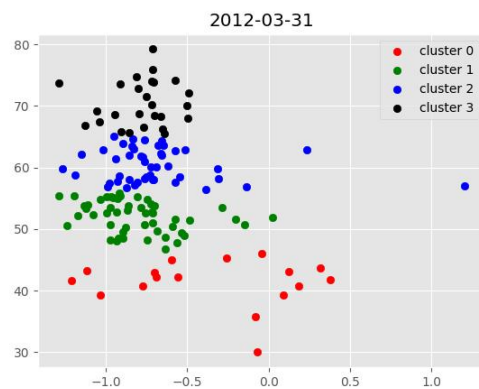
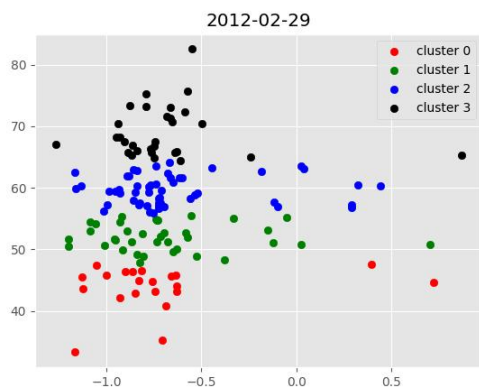
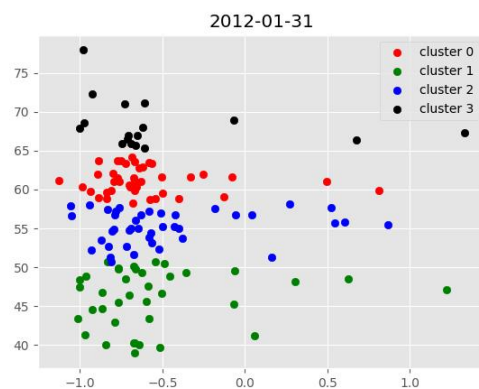
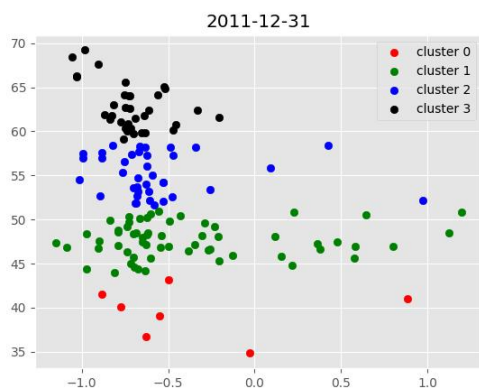
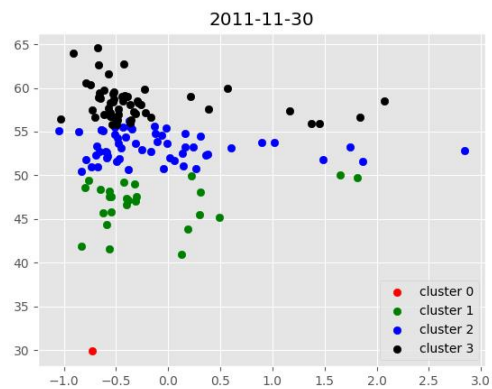
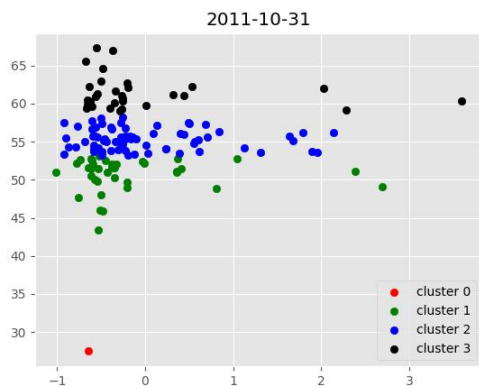
## **2.3 *K*-Means clustering and portfolio optimisation:**

We filter out stocks with  $< 10$  months of data, then calculate rolling factor betas and join it with the main features dataframe. Finally, we use an unsupervised *K*-Means clustering algorithm to group similar assets based on their features into 4 clusters. We have used the normalised values for all indicators except RSI, because we want to target specific RSI values. In fact our strategy is built on the hypothesis that momentum is persistent and that stocks clustered around RSI 70 centroid should continue to outperform in the following month - thus we select stocks corresponding to cluster 3.

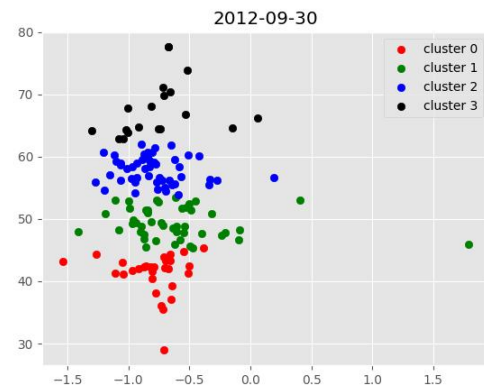
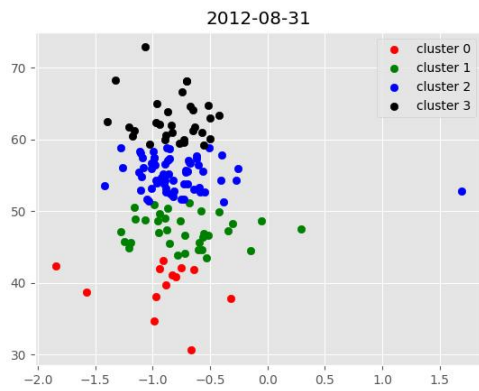
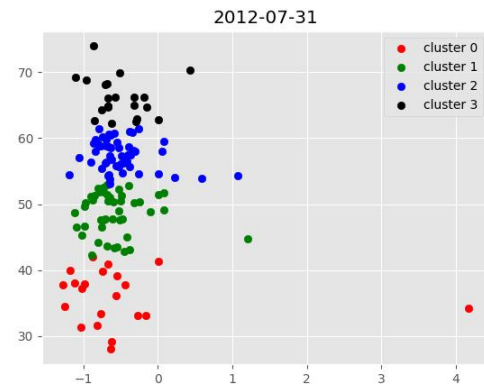
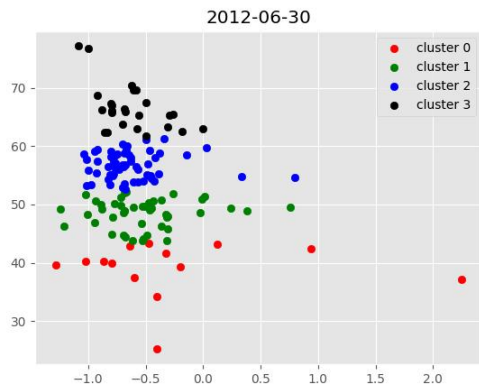
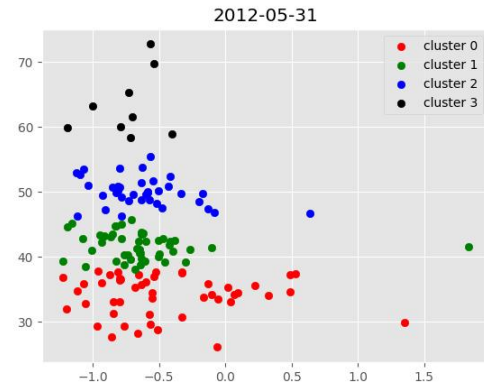
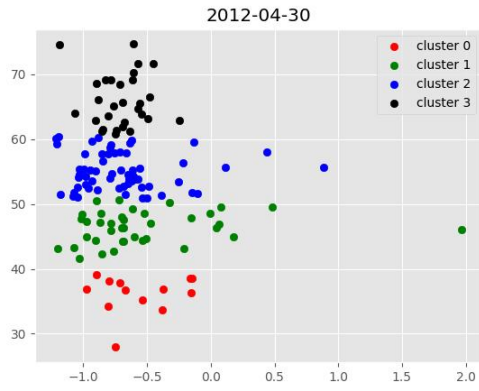
### 2.3.1 Cluster plots for each month

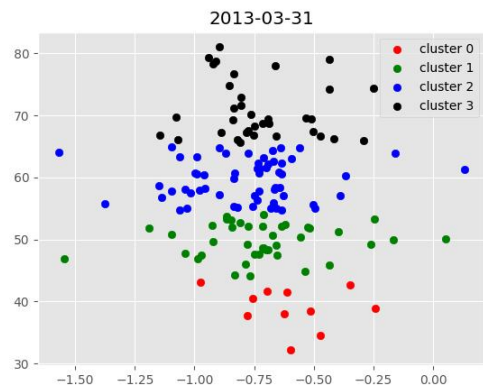
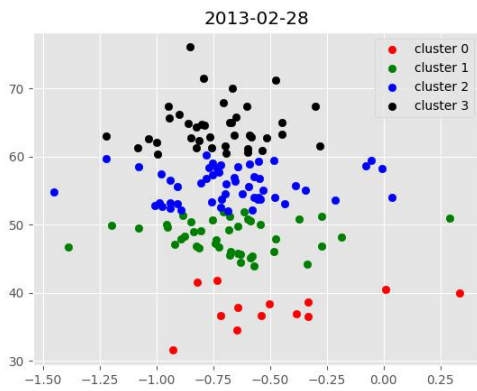
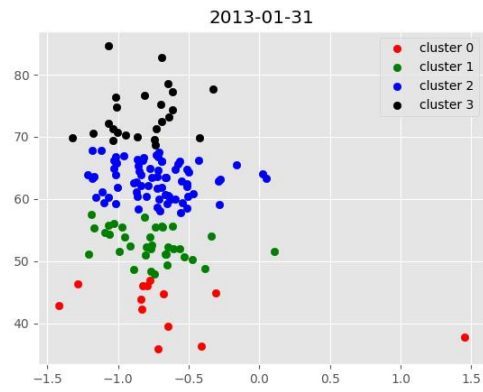
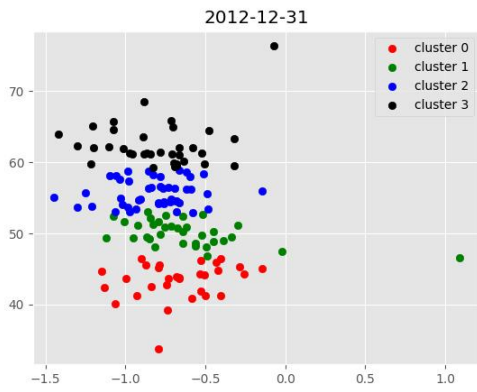
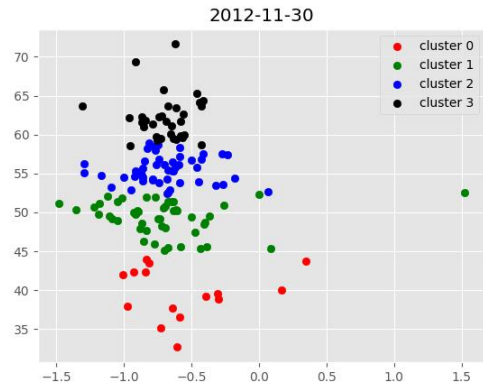
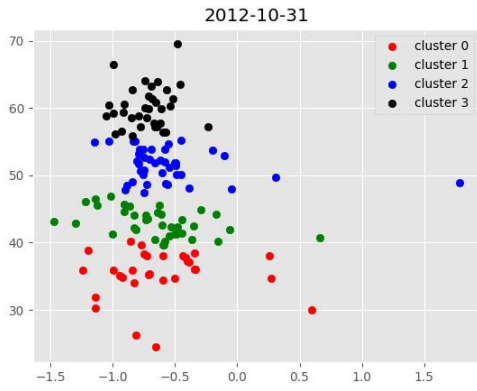


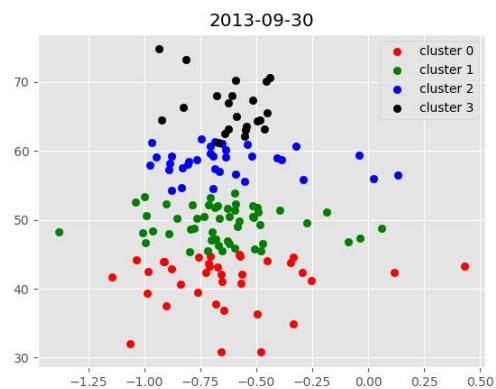
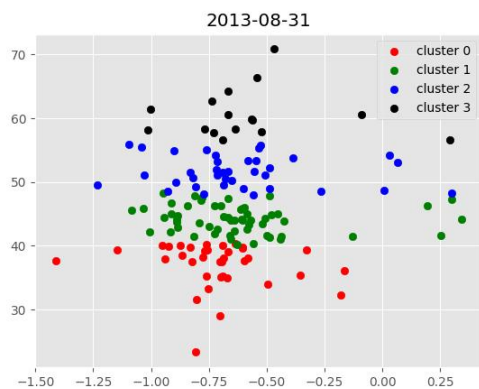
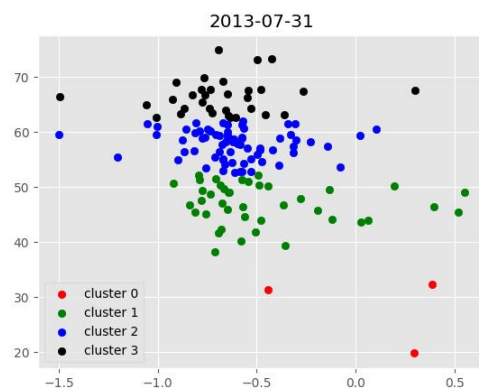
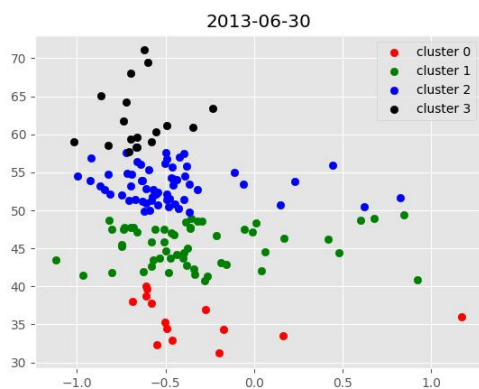
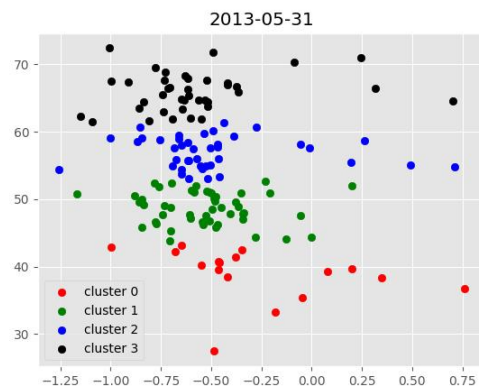
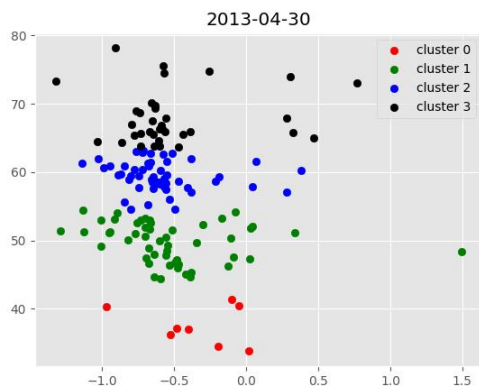


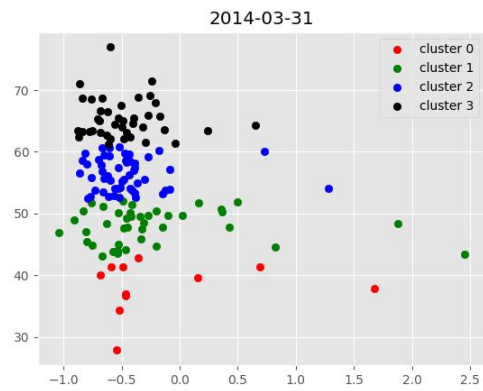
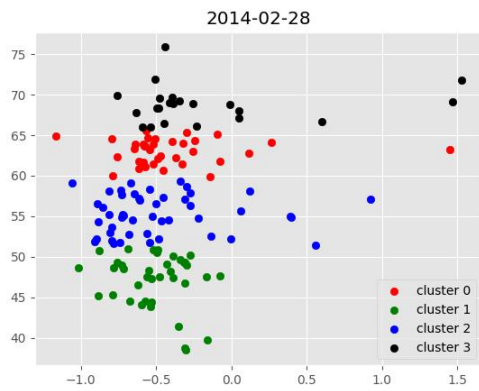
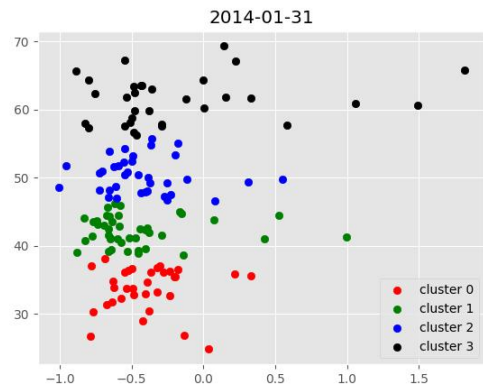
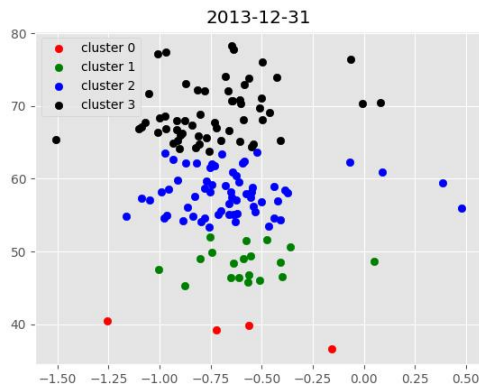
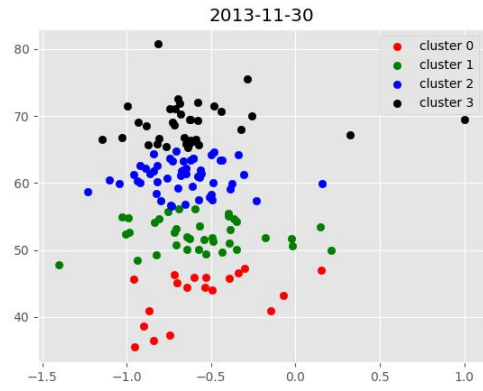
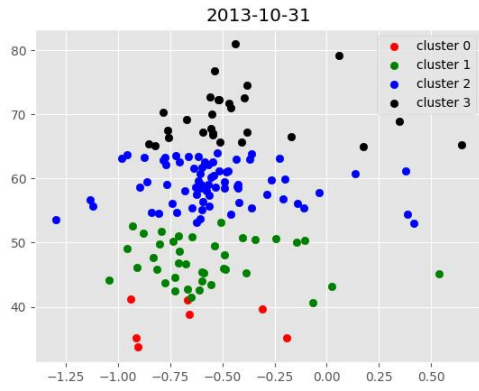


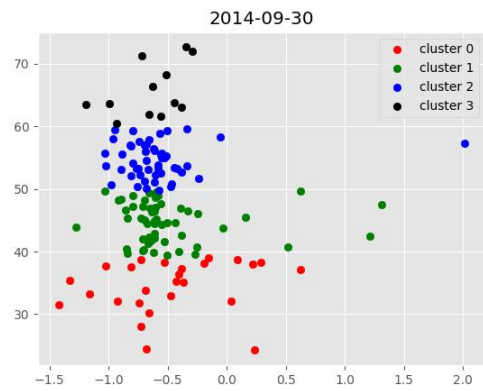
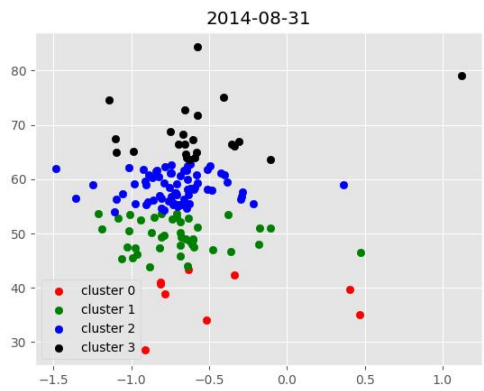
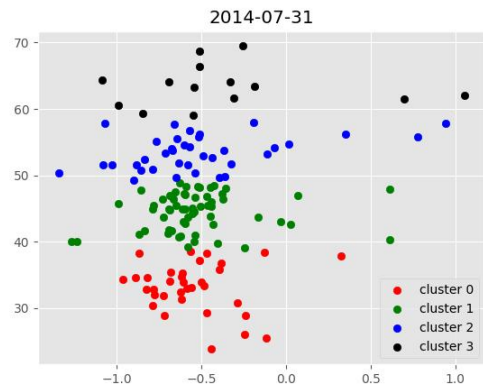
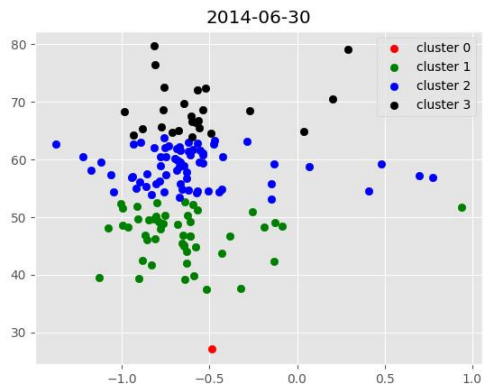
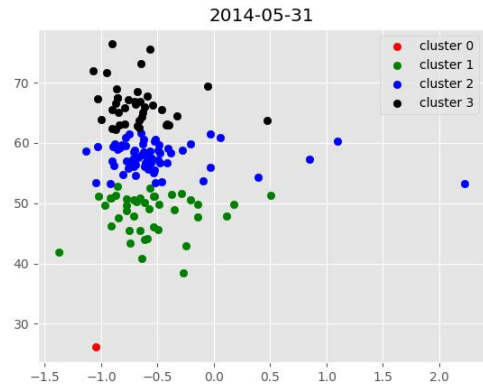
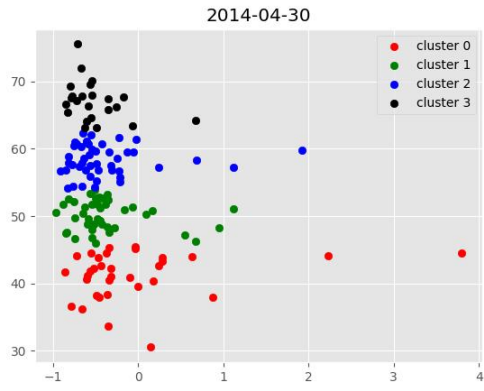


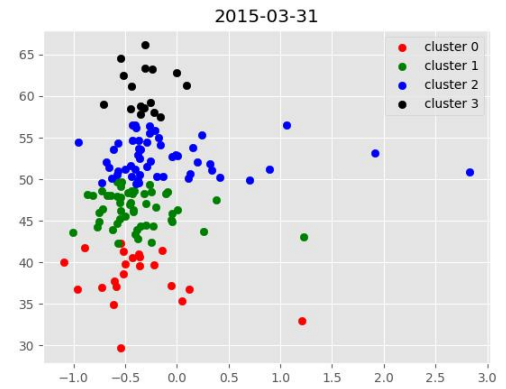
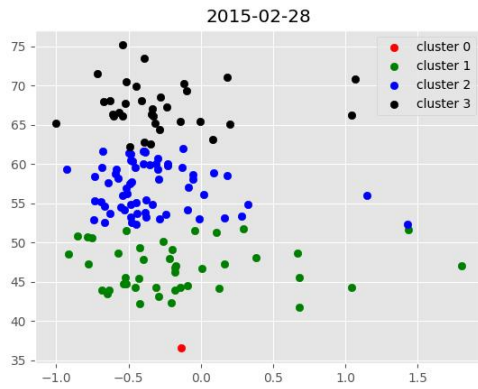
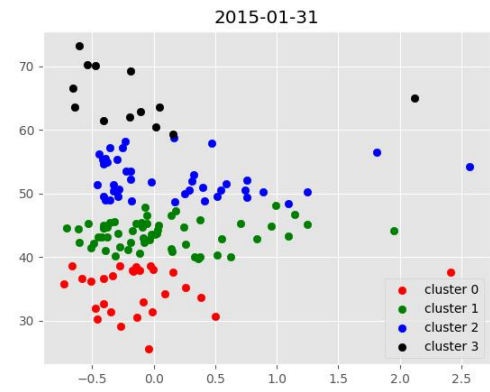
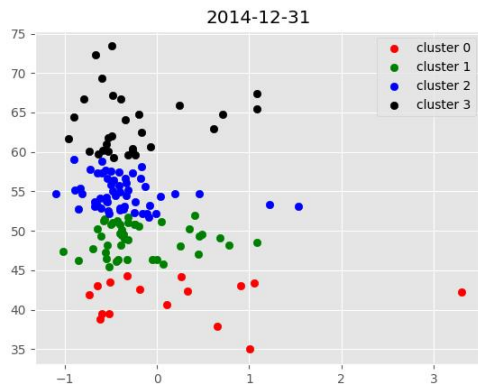
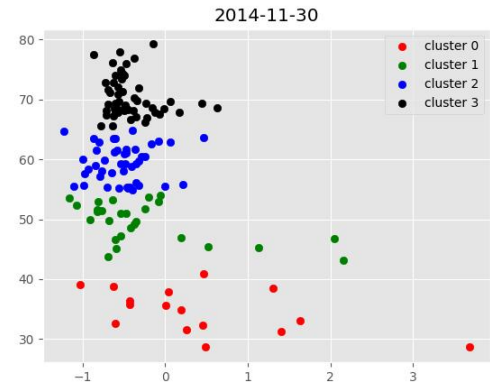
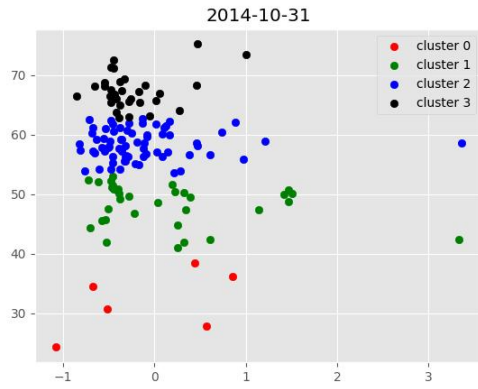


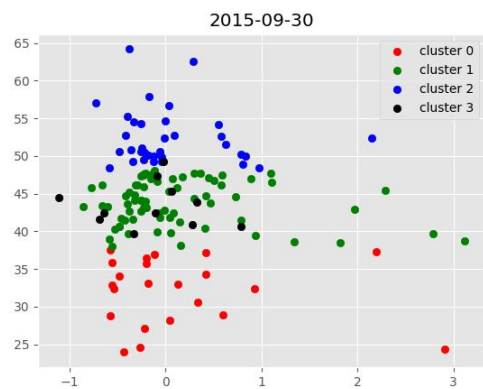
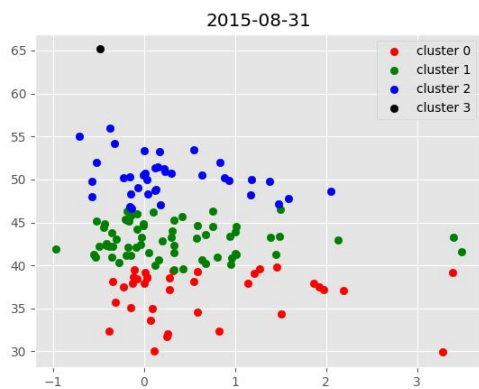
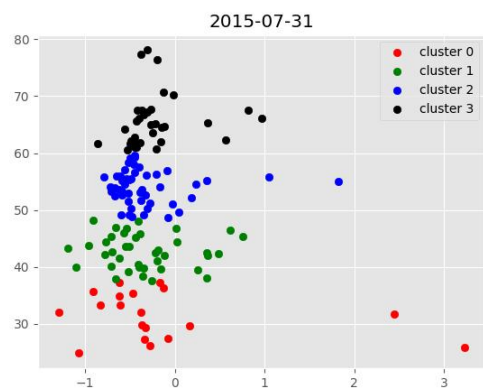
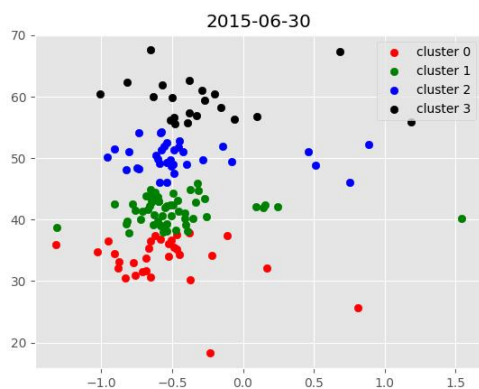
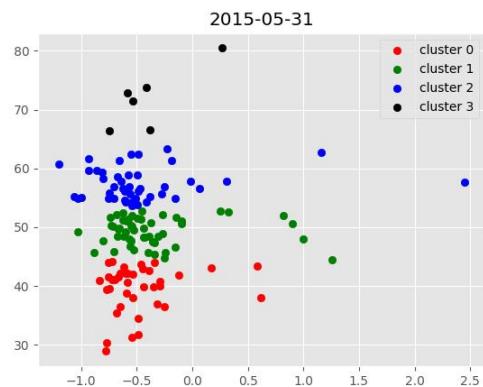
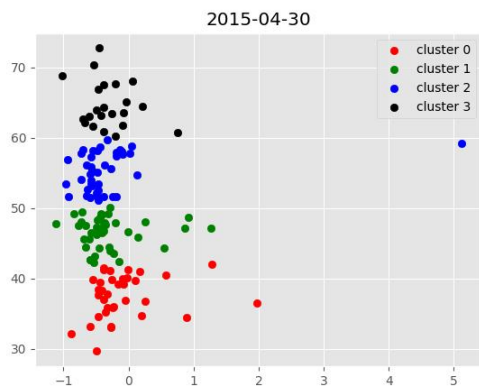




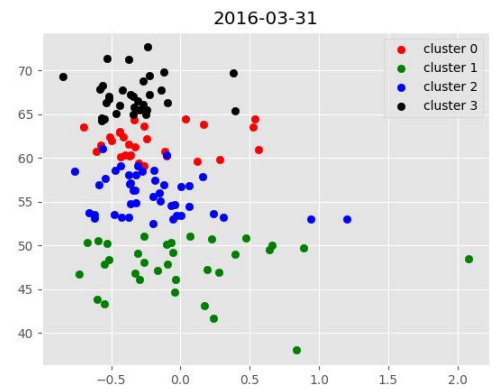
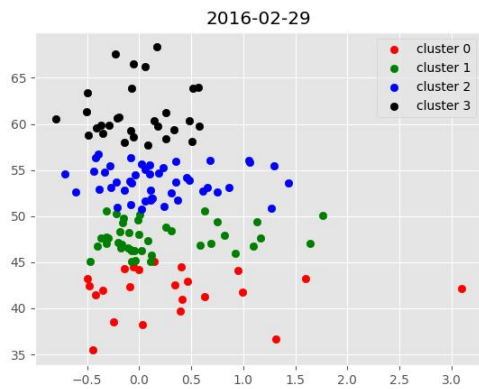
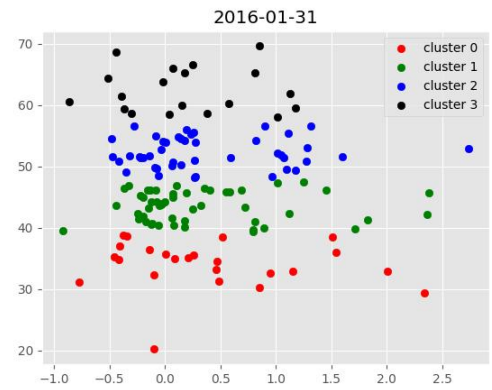
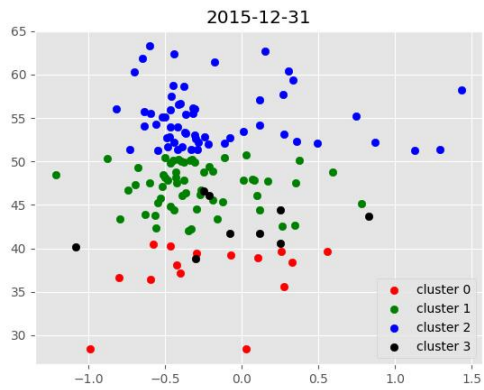
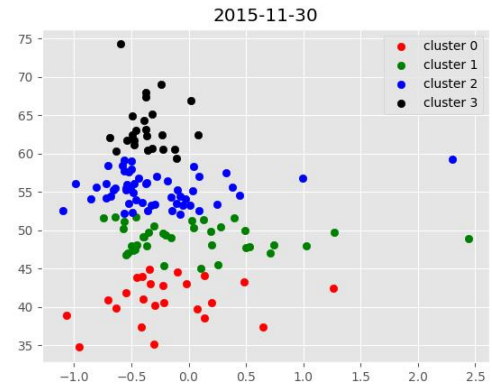
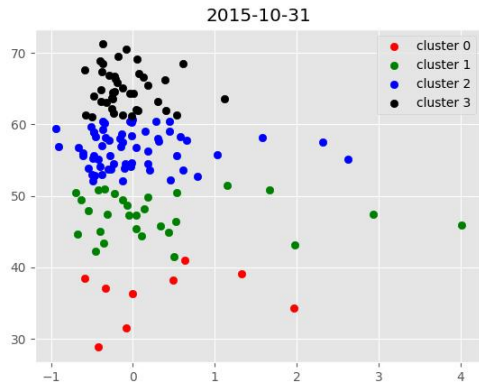




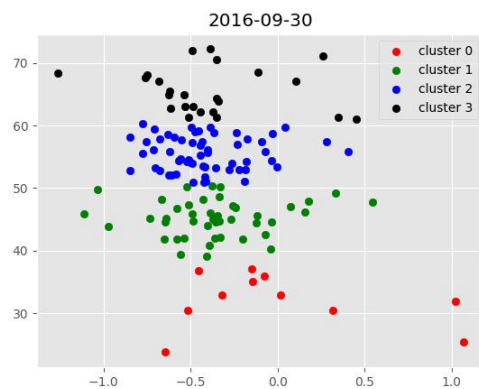
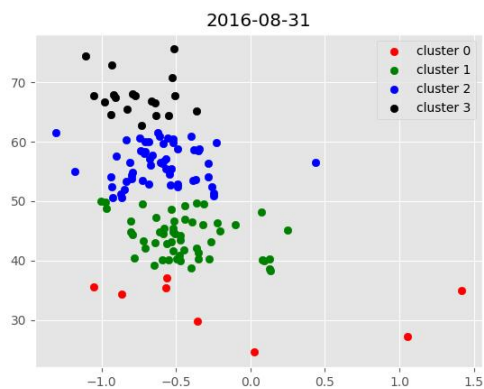
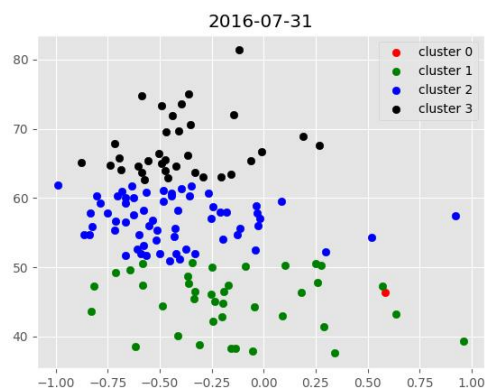
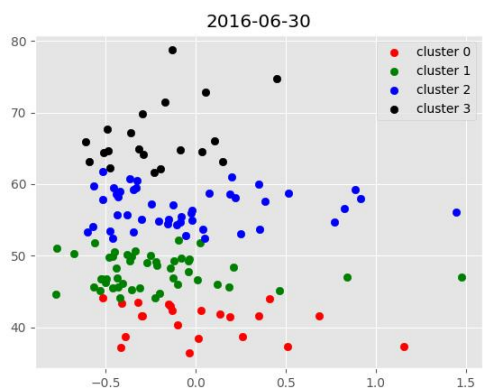
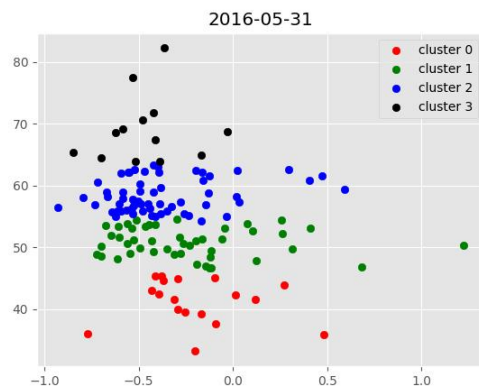
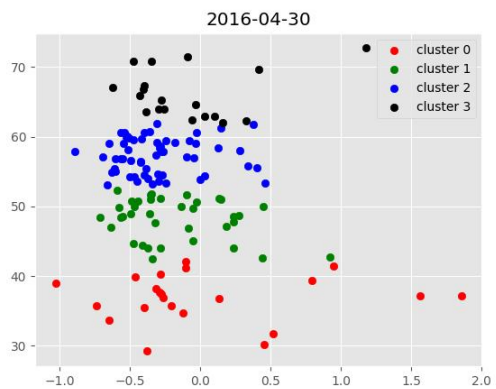


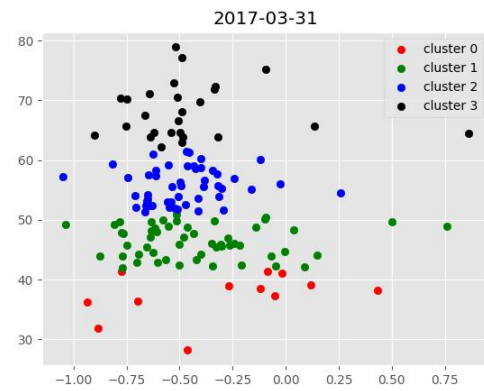
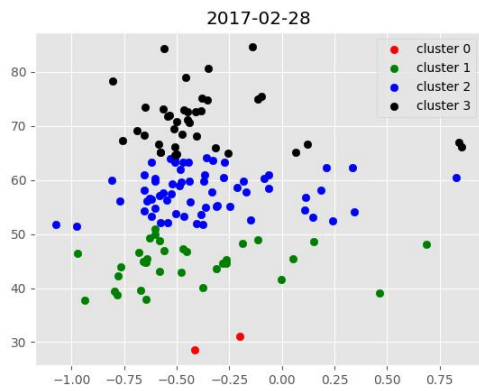
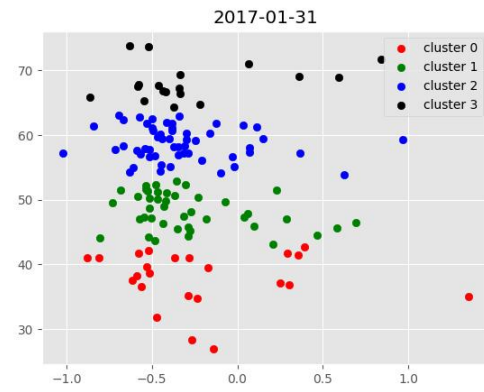
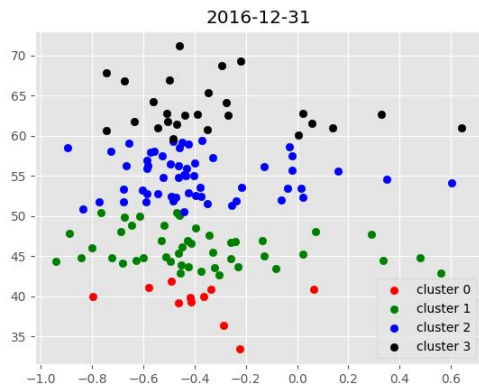
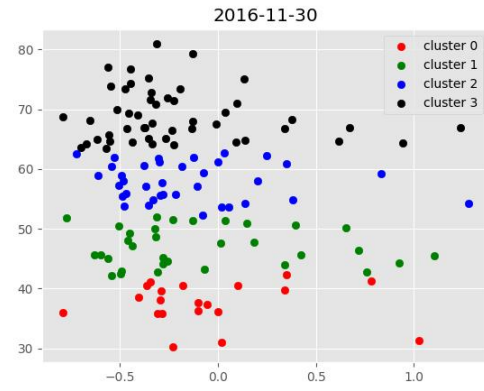
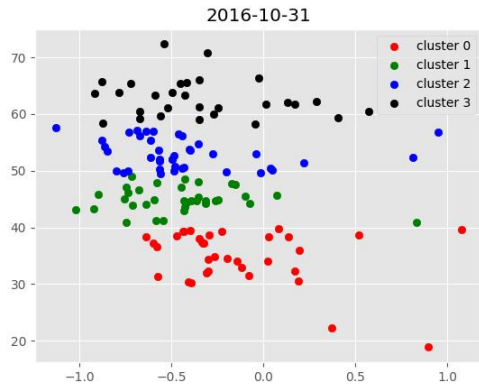


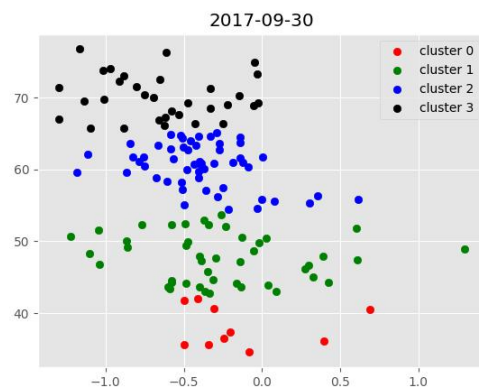
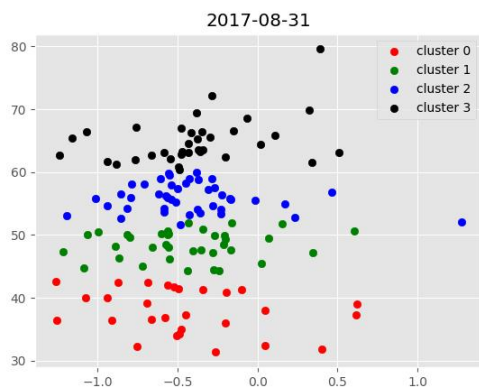
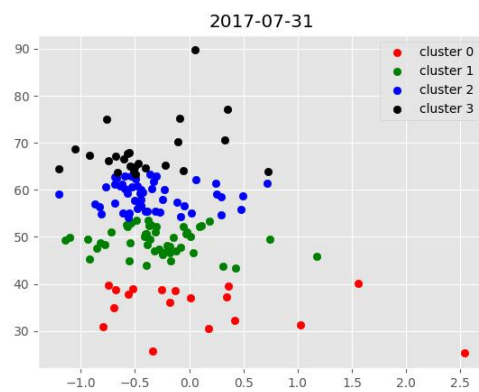
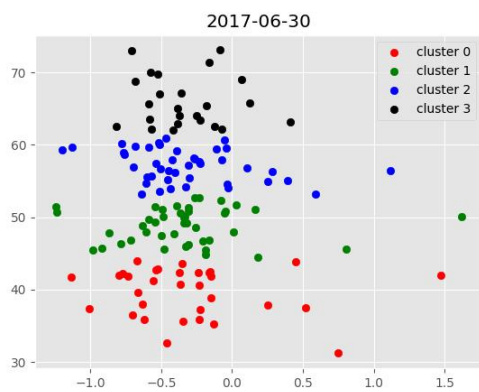
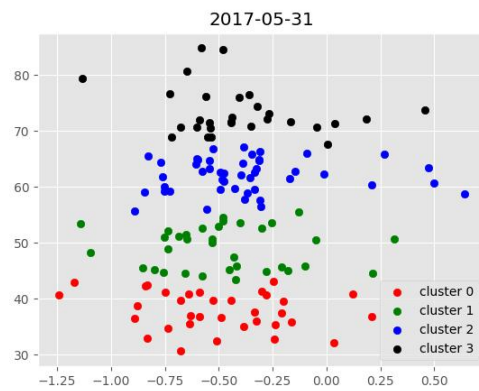
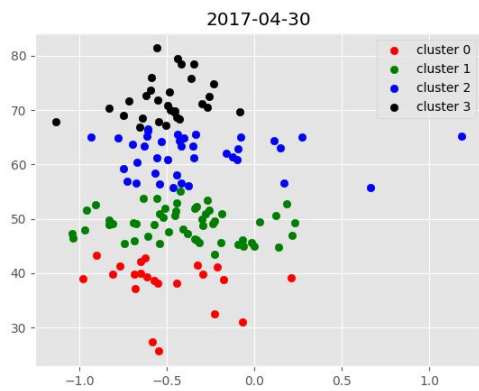


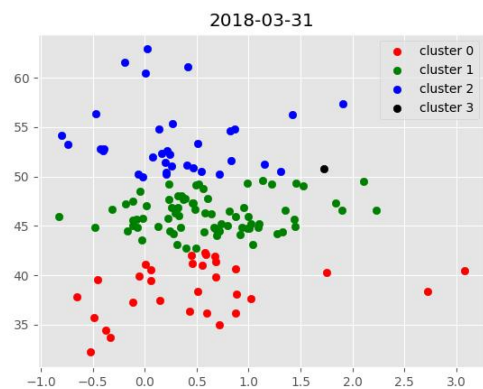
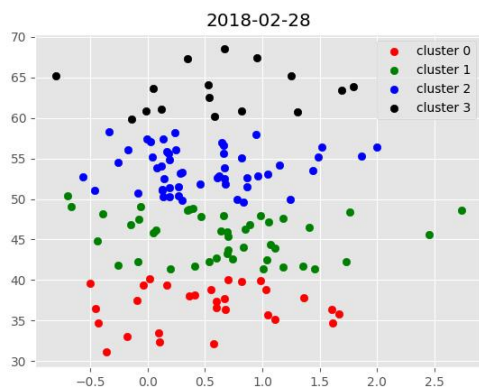
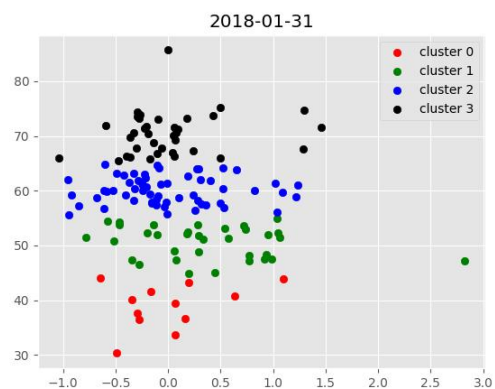
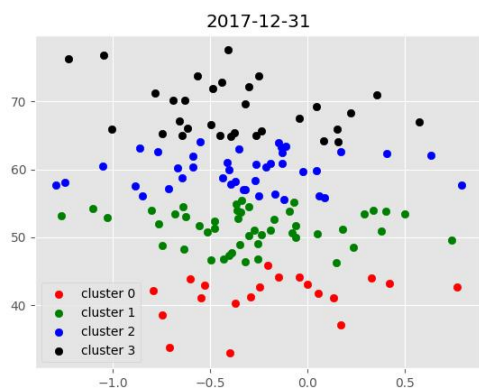
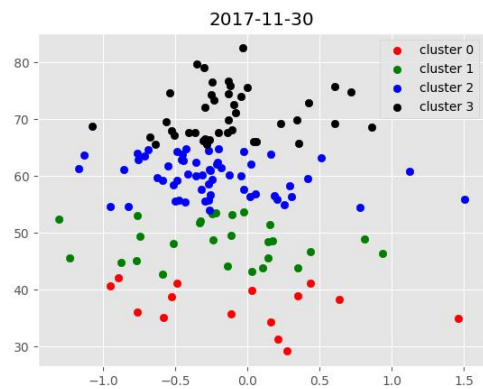
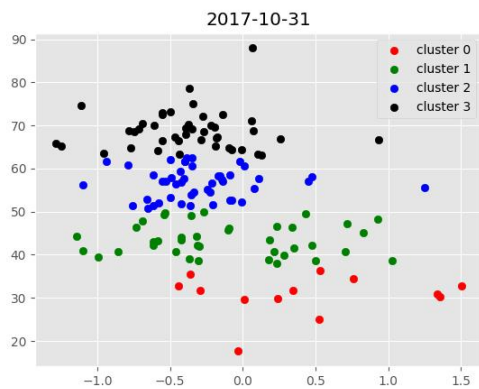


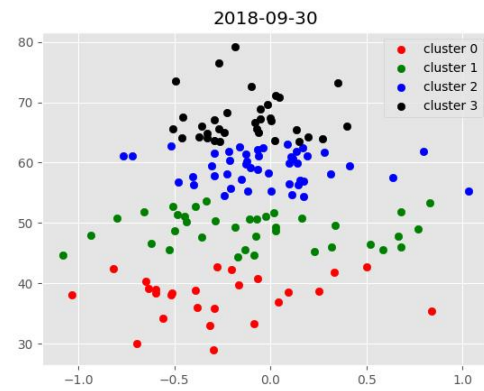
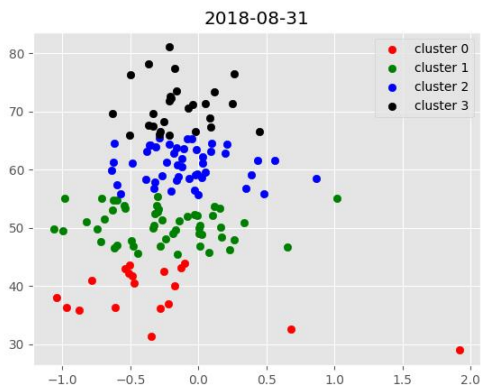
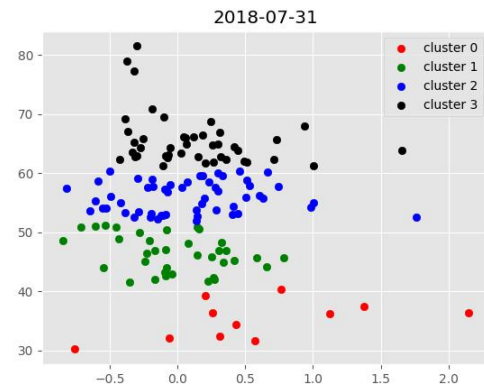
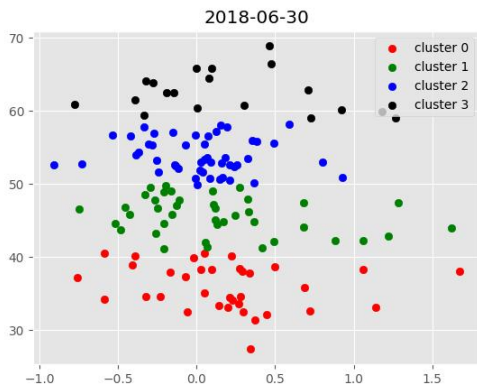
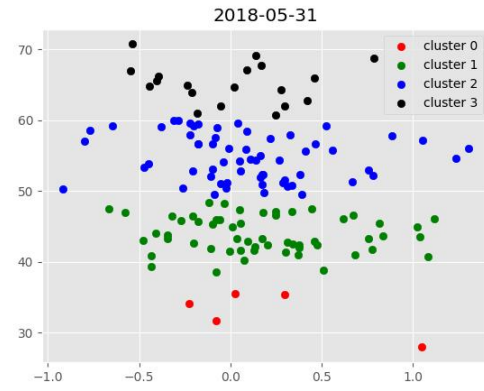
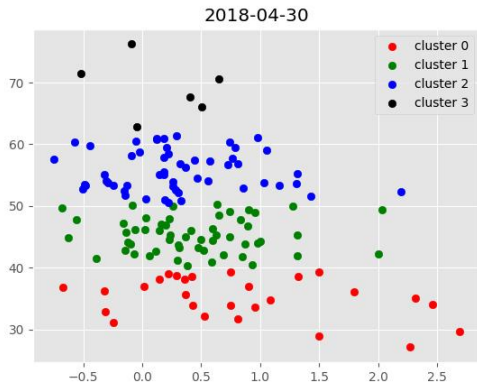


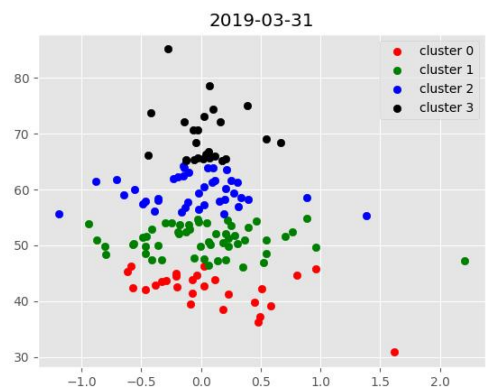
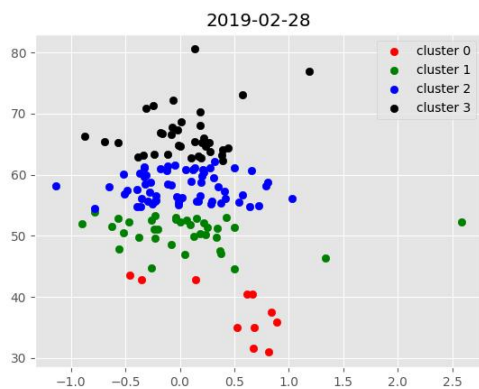
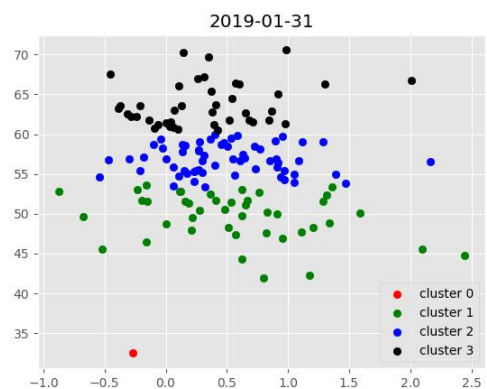
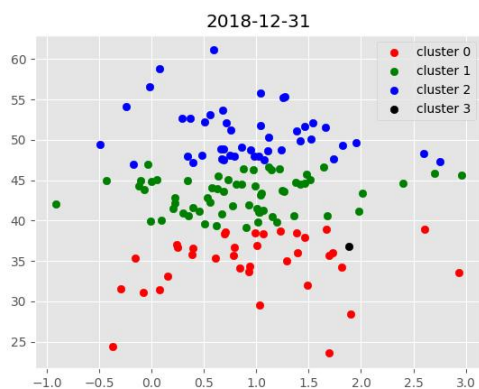
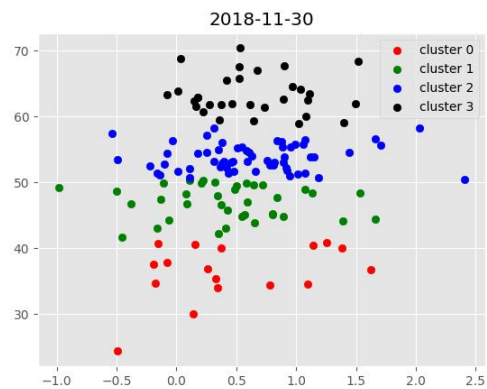
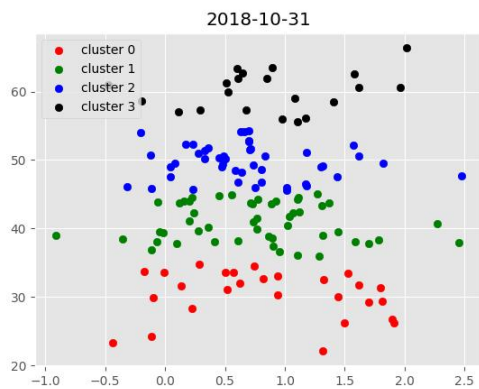




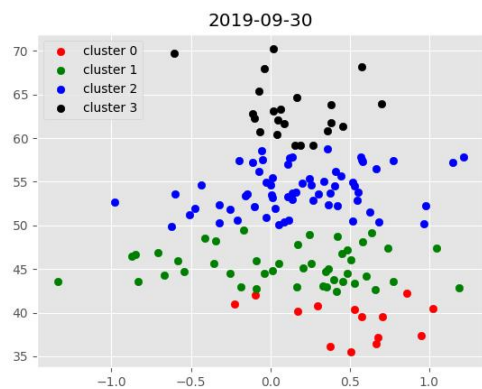
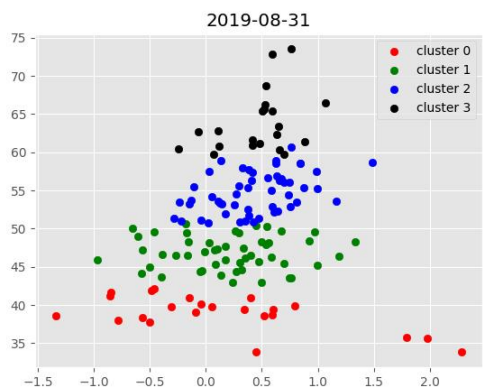
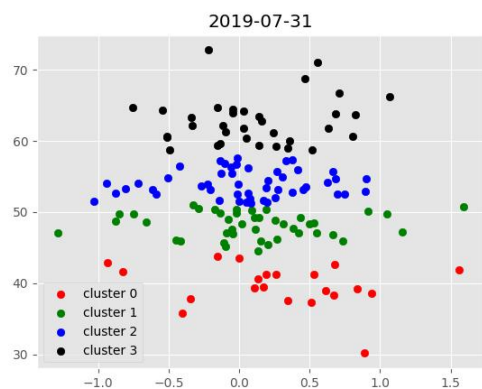
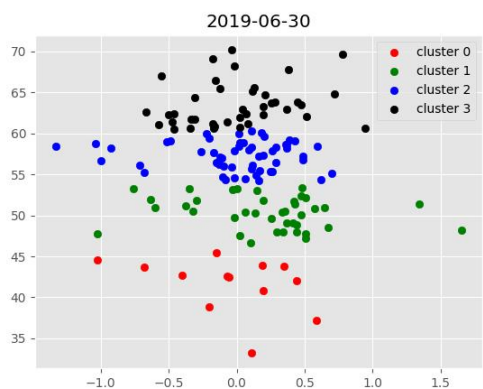
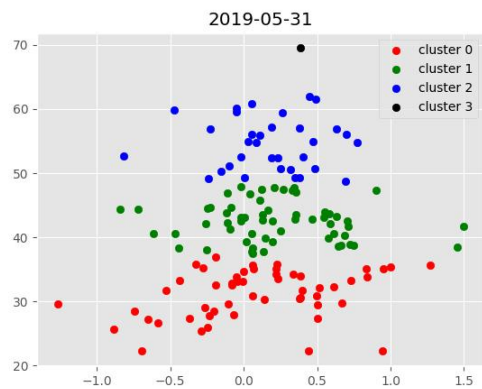
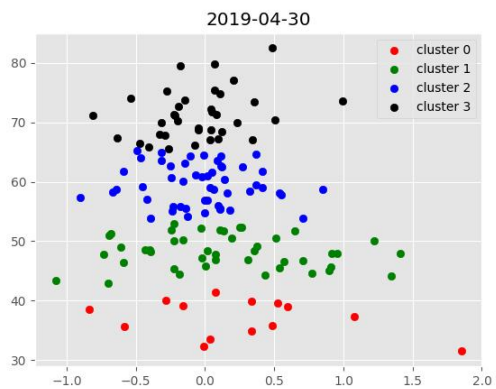


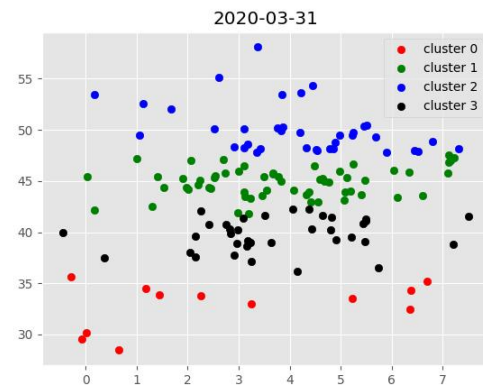
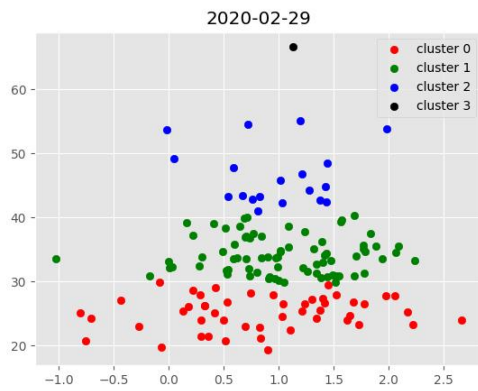
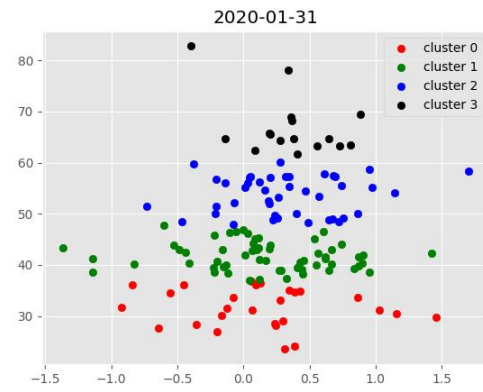
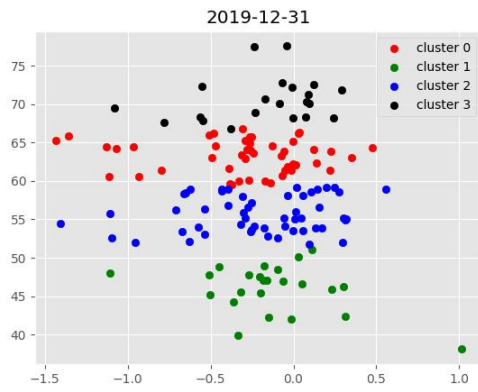
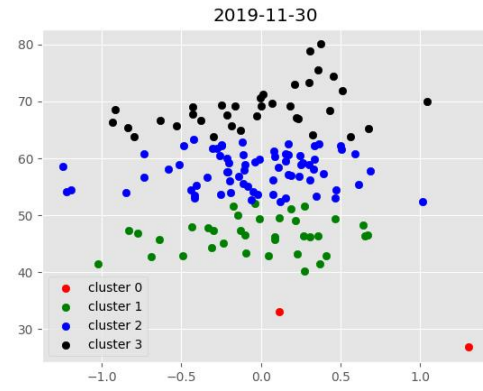
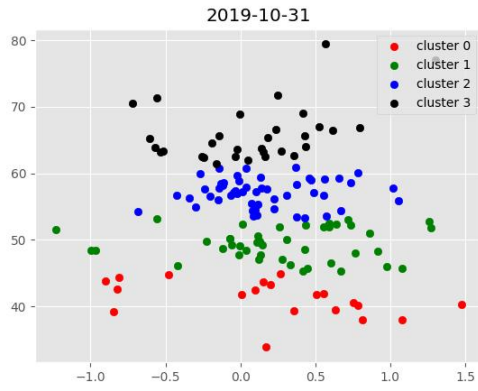




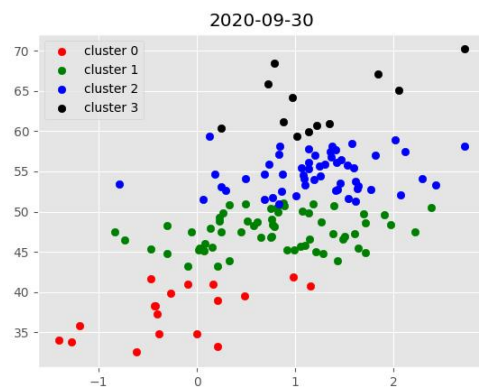
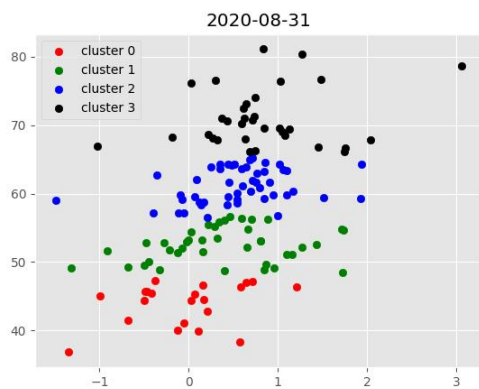
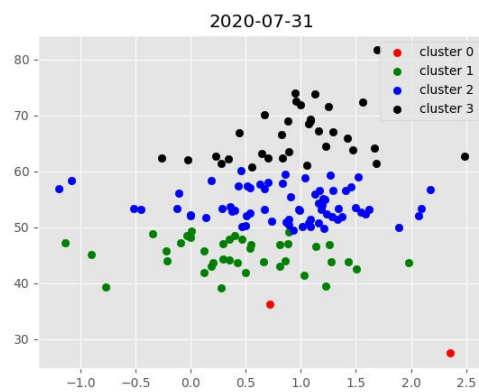
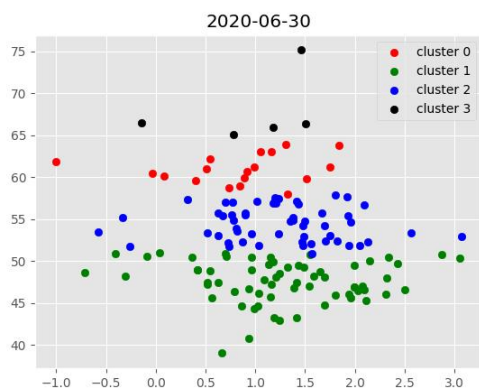
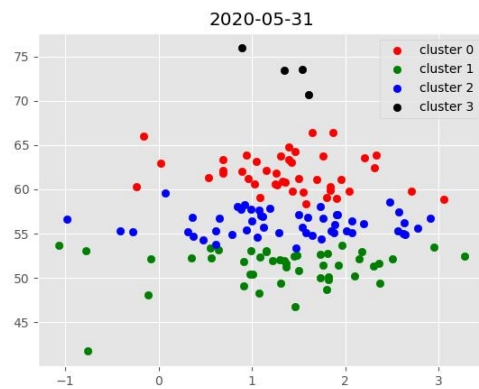
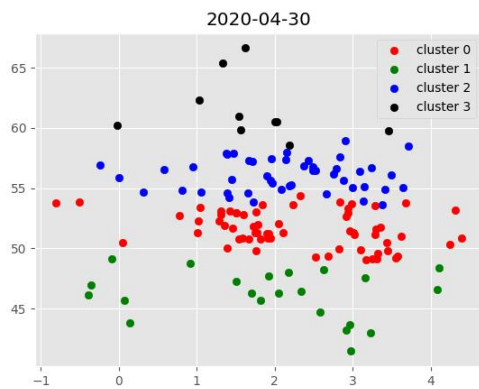


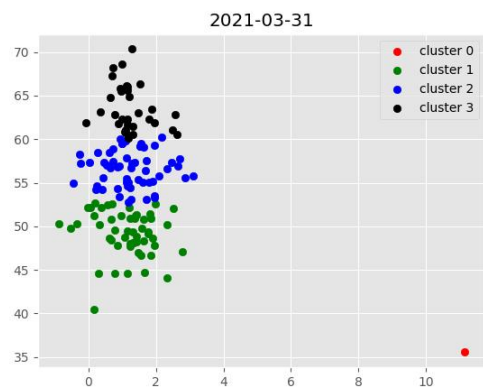
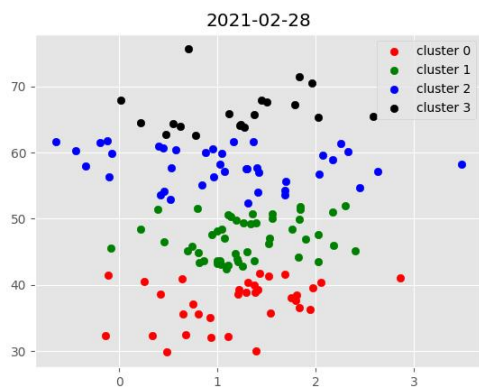
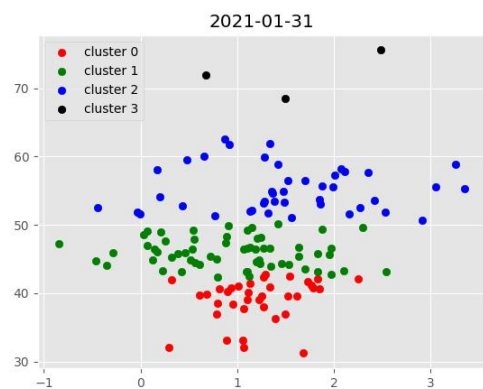
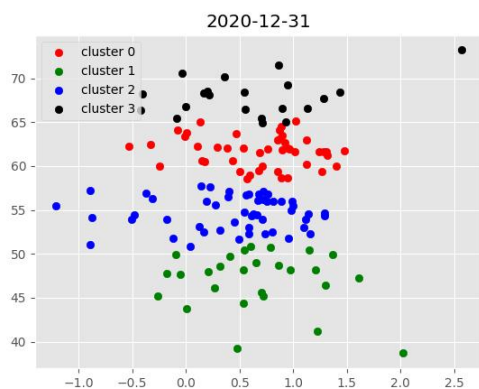
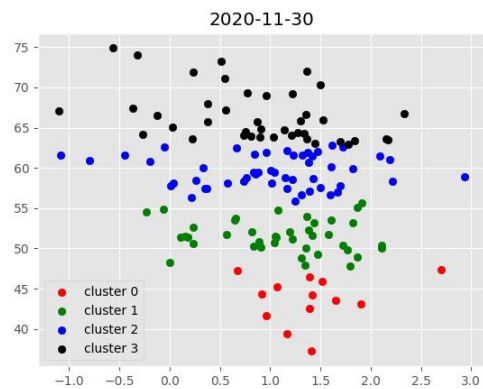
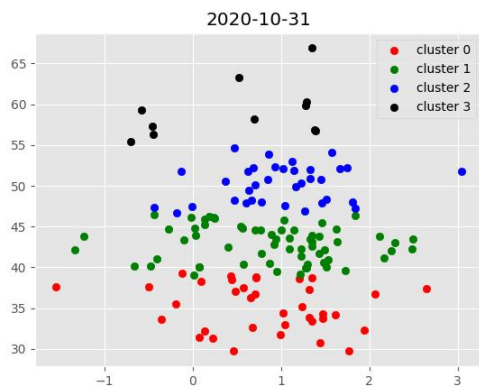


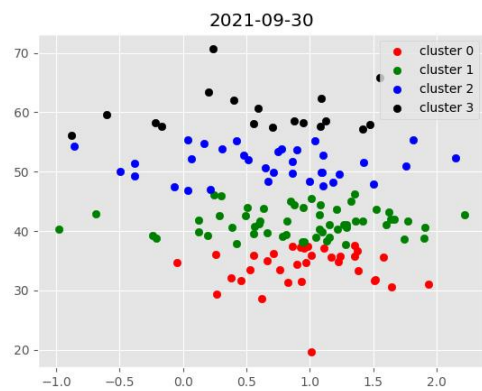
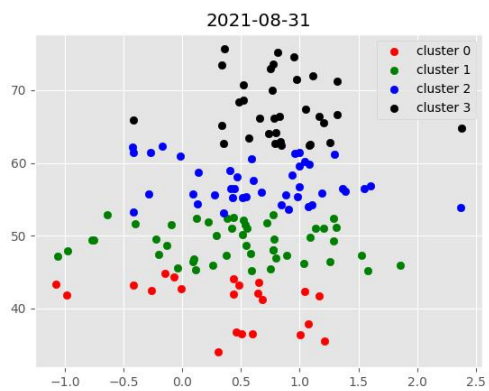
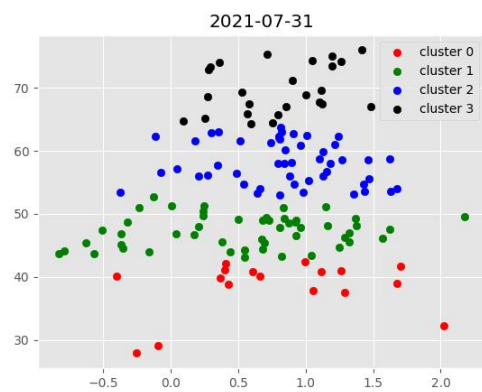
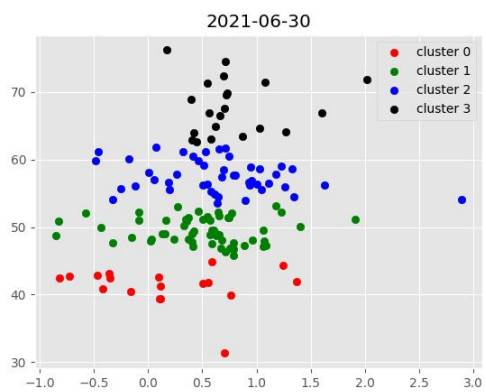
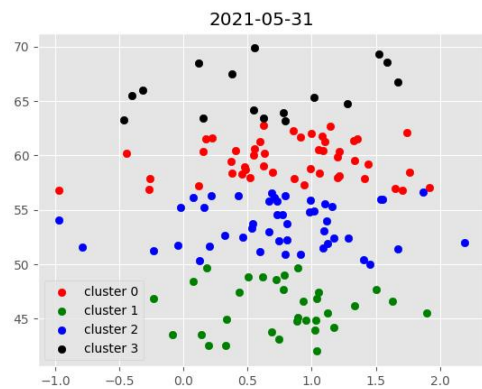
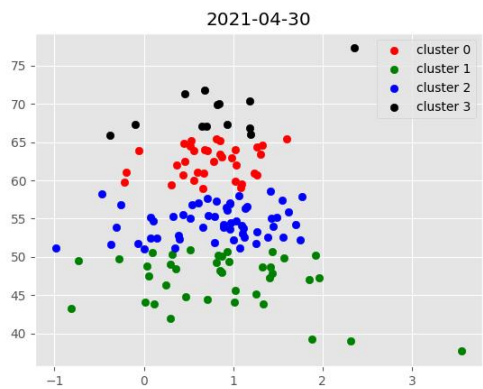


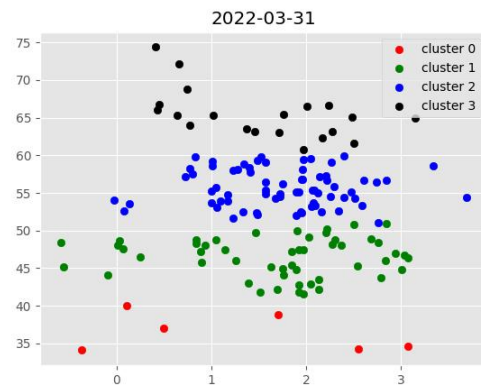
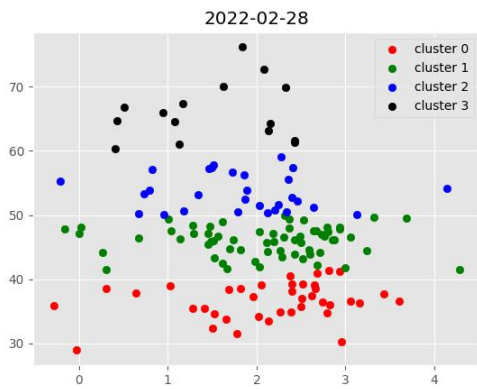
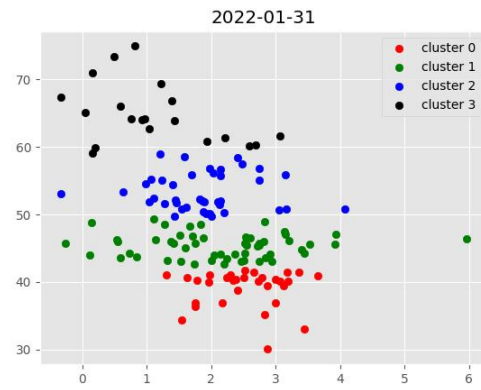
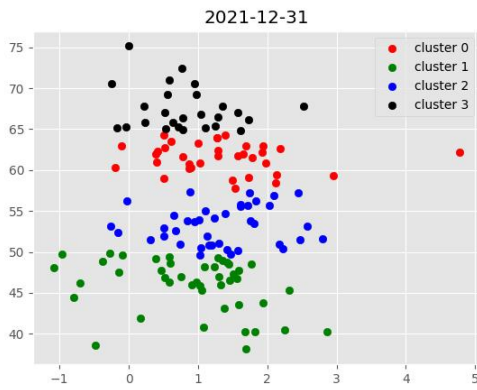
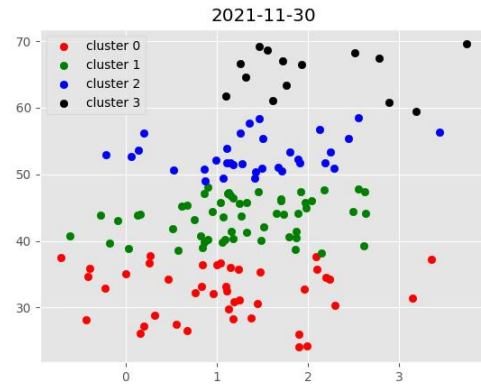
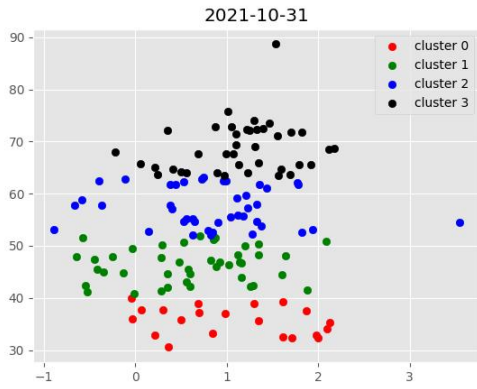


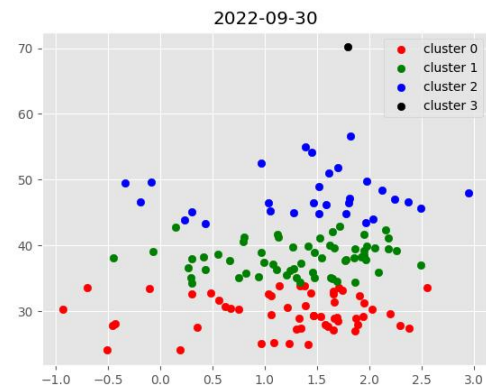
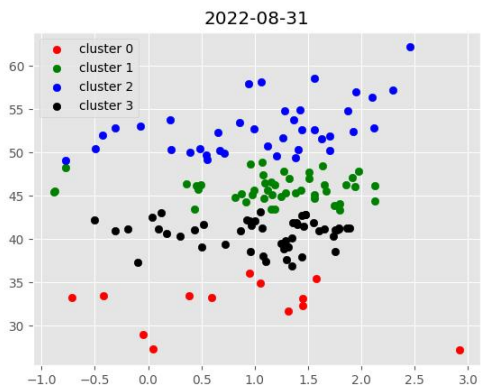
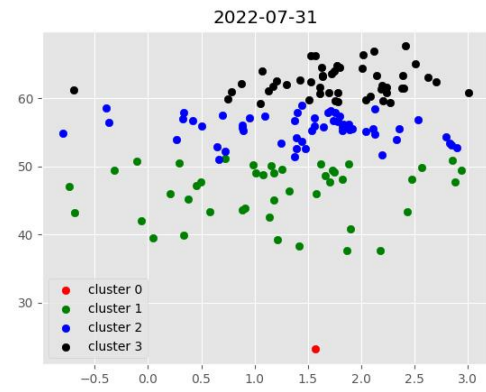
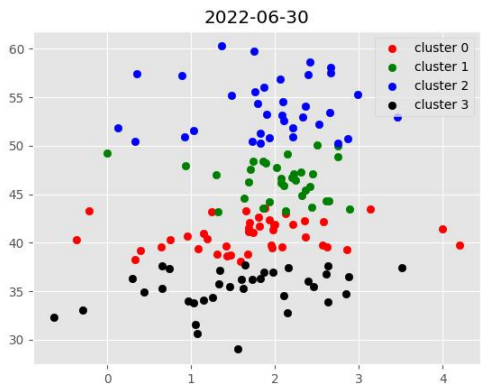
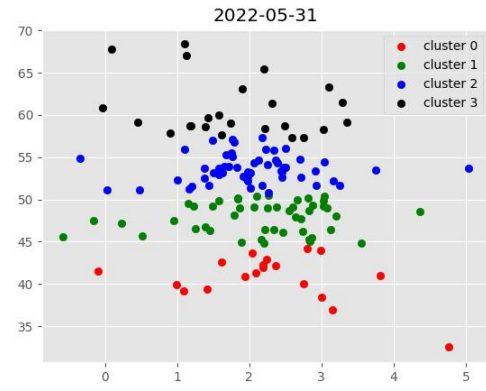
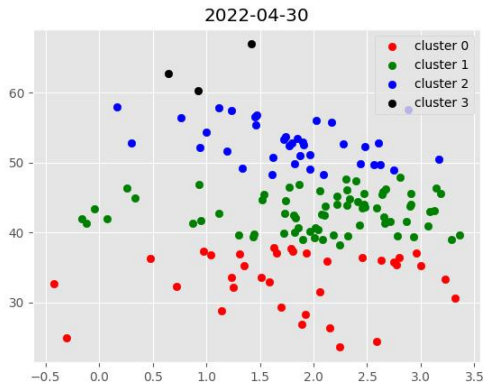


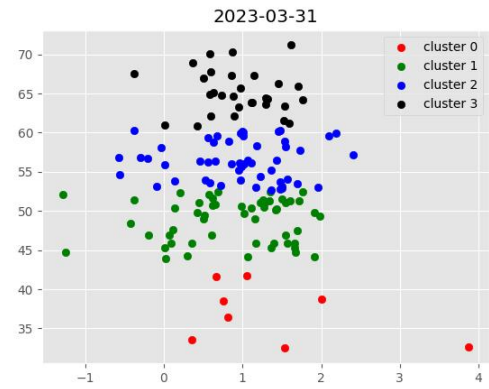
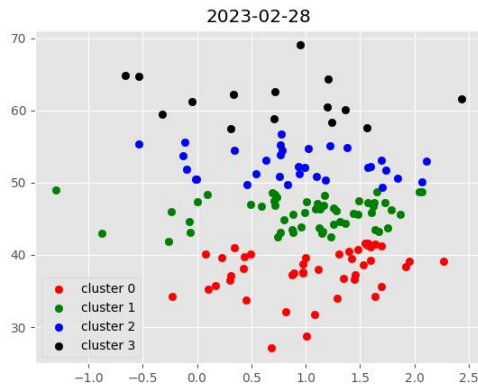
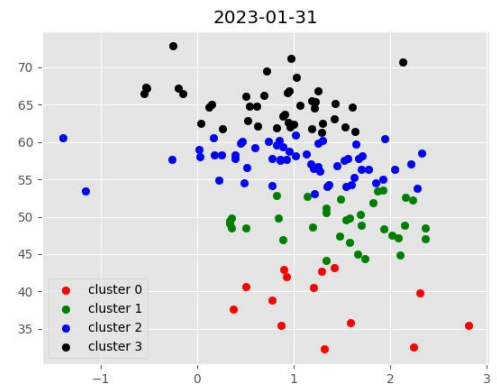
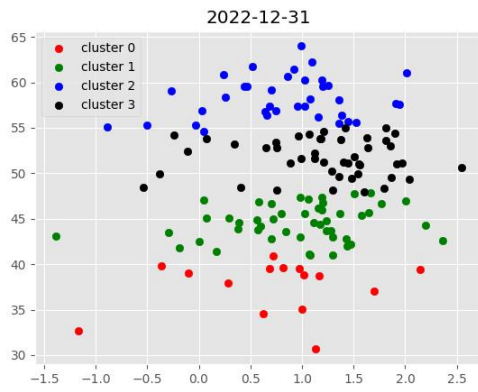
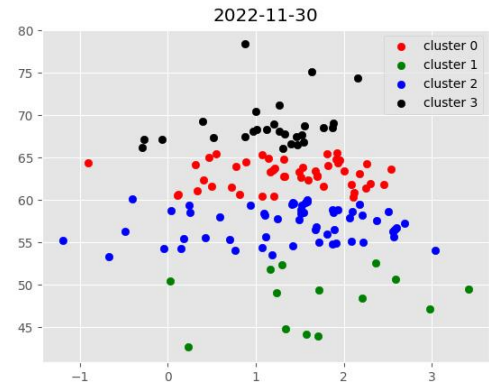
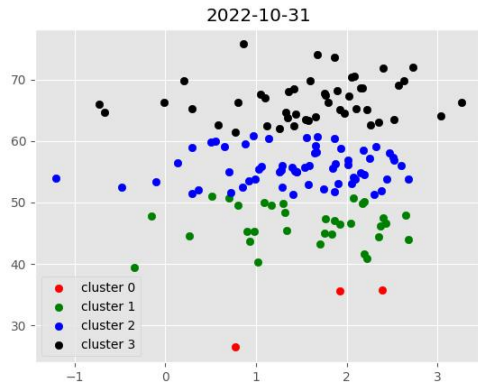




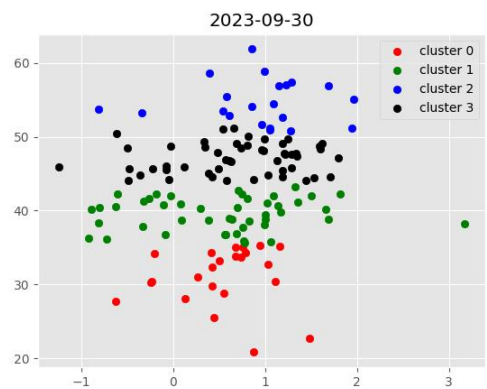
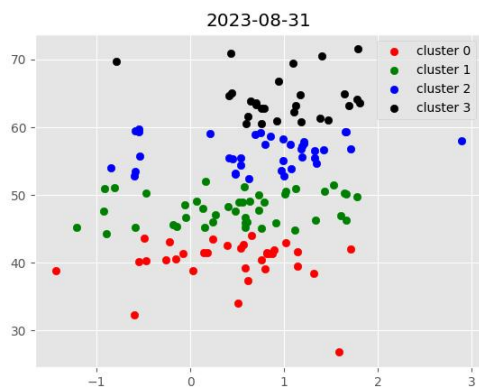
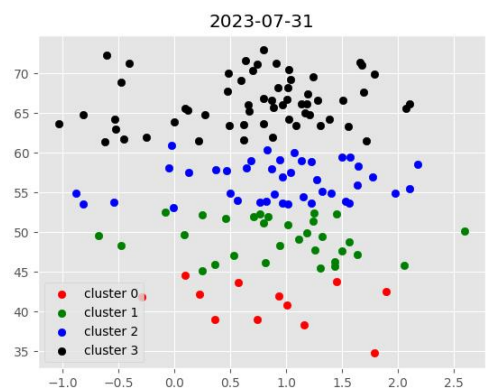
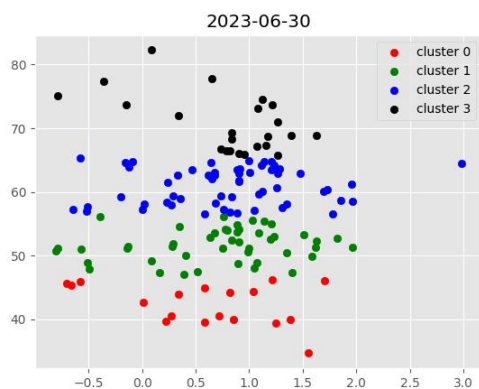
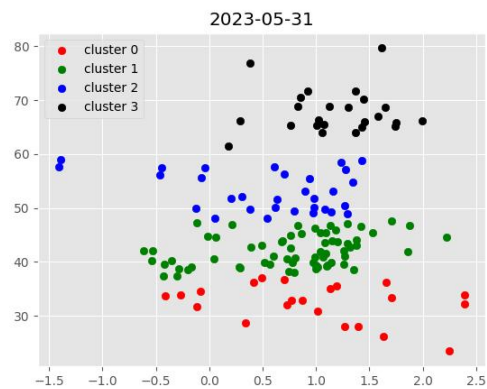
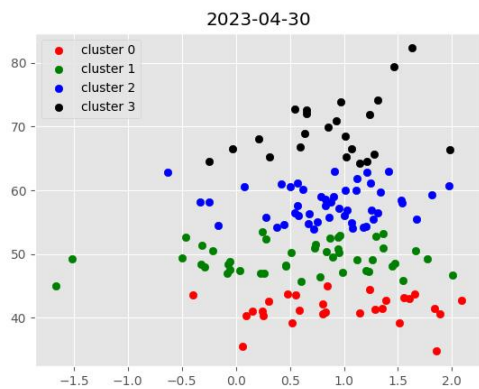


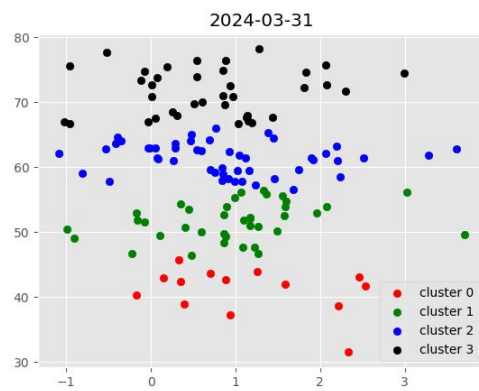
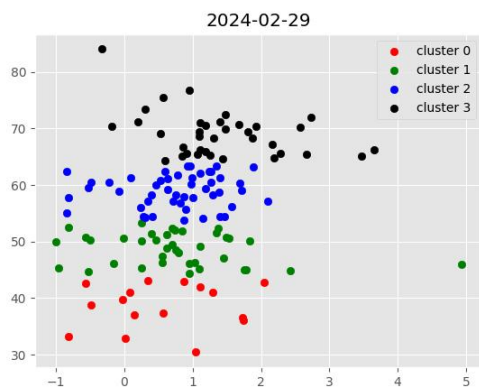
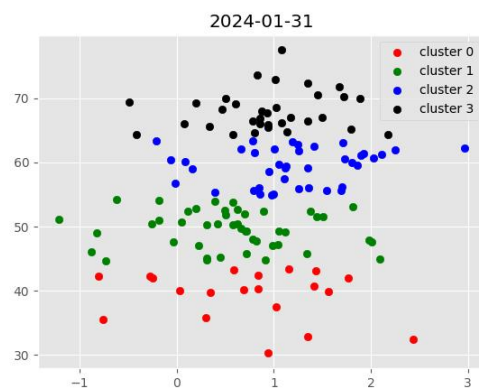
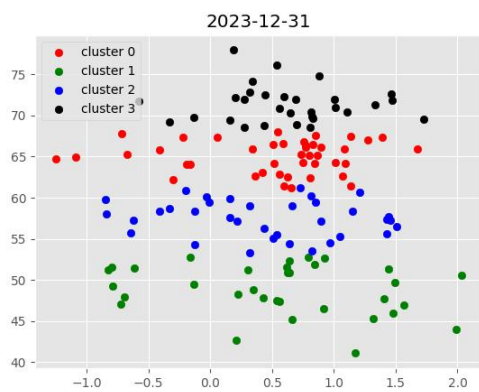
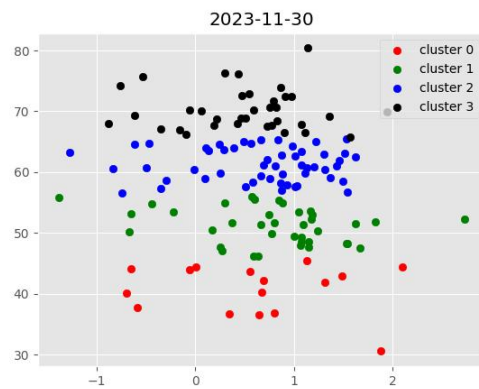
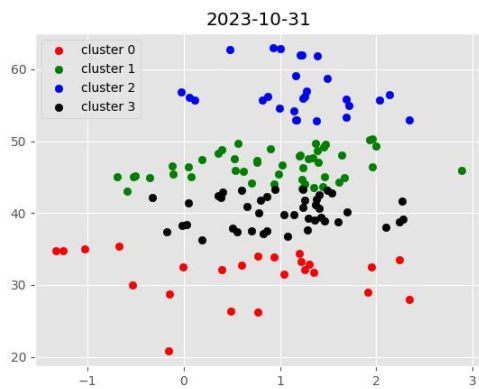




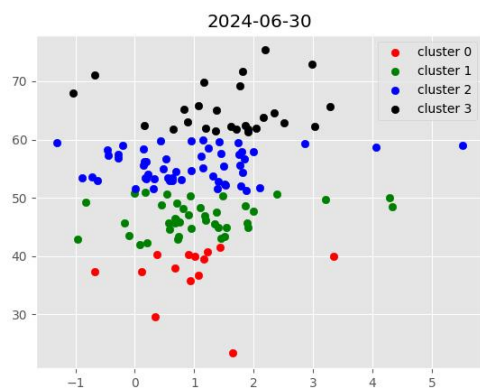
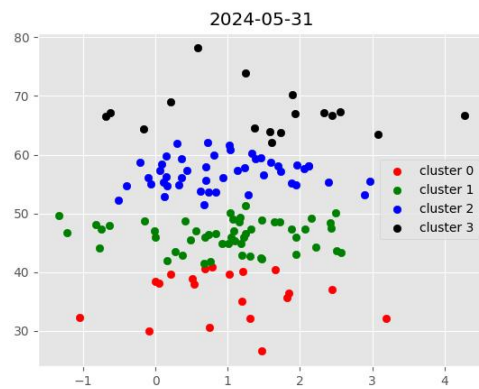
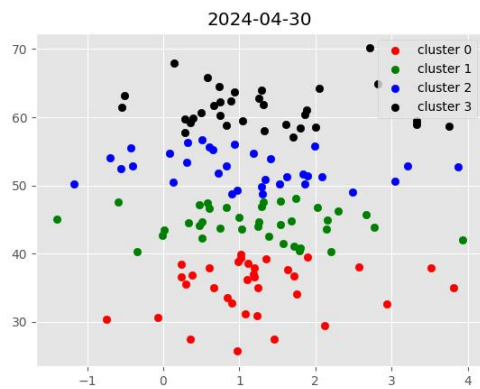














For each month we select assets based on the cluster and form a portfolio based on efficient frontier maximum Sharpe ratio optimisation.

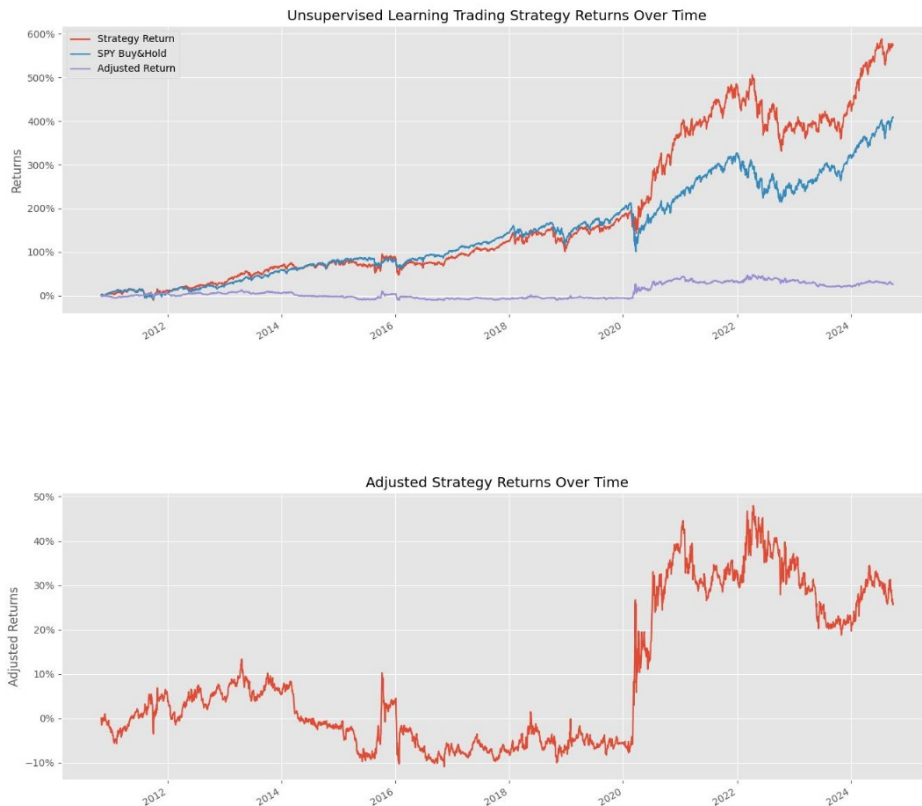
- First we will filter only stocks corresponding to the cluster we choose based on our hypothesis.
- As per our strategy hypothesis we select stocks corresponding to cluster 3.
- We will define a function which optimises portfolio weights using PyPortfolioOpt package and EfficientFrontier optimiser to maximise the Sharpe ratio.
- To optimise the weights of a given portfolio we would need to supply last year's prices to the function.
- Apply single stock weight bounds constraint for diversification (minimum half of equally weighted and maximum 10% of portfolio).

## **2.4 Computing portfolio returns:**

- Calculate daily returns for each stock which could land up in our portfolio.
- Then loop over each month's start, select the stocks for the month and calculate their weights for the next month.
- If the maximum Sharpe ratio optimization fails for a given month, apply equally-weighted weights.
- Calculate daily portfolio returns.

### 3 Final results

Finally we visualise portfolio returns and compare to S&P 500 returns.



#### 3.1 Limitations and possible future extensions

We do not consider stocks that have been delisted, so there is survivorship bias in our dataset. Yahoo Finance API has been officially discontinued and the current package that we use (yfinance) is an open-source tool that uses Yahoo's publicly available APIs. It is not perfect and returns N/A values for some crucial quantities which is why we had to drop N/A values from the dataframes in our code at various points. This could be mitigated by using a better API with high quality data, but almost all of them are pay-to-use. Also, we haven't split our data so there could be some overfitting that we haven't accounted for.

The model fails to outperform the SPY Buy&Hold over shorter timescales (such as  $< 5$  years). This could be mitigated by extending our current model to use neural networks or other models, and picking the most optimal strategies from the ones given by each model.

## 4 References

- [1] M. F. Dixon, I. Halperin, and P. Bilokon, Machine Learning in Finance. Springer International Publishing, 2020. doi: 10.1007/978-3-030-41068-1.
- [2] FreeCodeCamp. Algorithmic Trading – Machine Learning & Quant Strategies Course with Python. (Oct. 26, 2023). Accessed: Sep. 25, 2024. [Online video]. Available: <https://www.youtube.com/watch?v=9Y3yaoi9rUQ>

