

# candlestick\_chart\_recognition

May 17, 2019

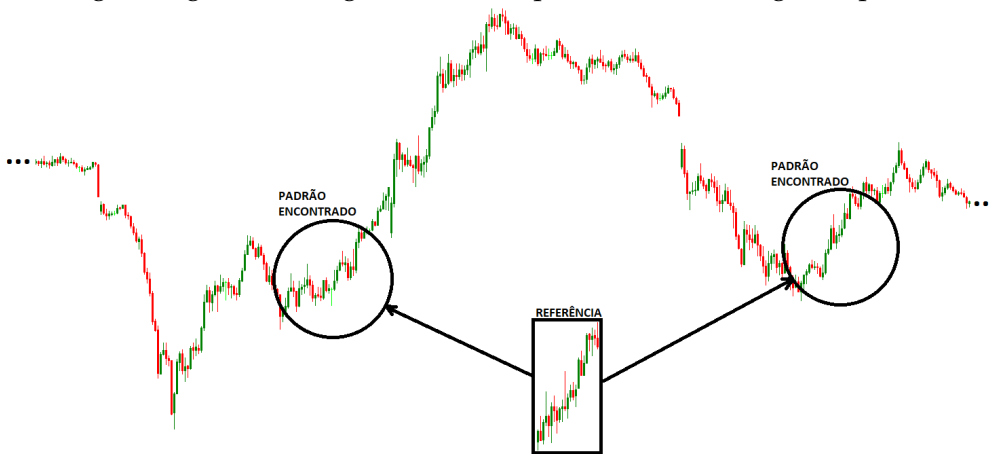
## 1 Projeto 'Candlestick Set Recognition'

1.0.1 Erivelton Marques de Carvalho

17-05-2019

### 1.1 Resumo

Este projeto executa um processo para reconhecer um gráfico pequeno, sendo este a referência, dentro de outro bem maior. Podem ser encontradas em várias partes. O tipo de gráfico que iremos trabalhar é chamado de 'candlestick', aplicado no mercado financeiro, especificamente nas bolsas de valores. O sistema executará uma varredura, ponto a ponto, do início ao fim do gráfico maior e armazenará o conjunto de pontos que tem alta correlação com os pontos do gráfico de referência. Abaixo, segue a figura de um gráfico desse tipo com uma analogia do processo.



A partir dos resultados, o investidor pode tomar algumas decisões, como comprar um determinado lote de ações de uma empresa considerando que, segundo os padrões encontrados no gráfico de histórico (grafico maior), o preço tende a subir.

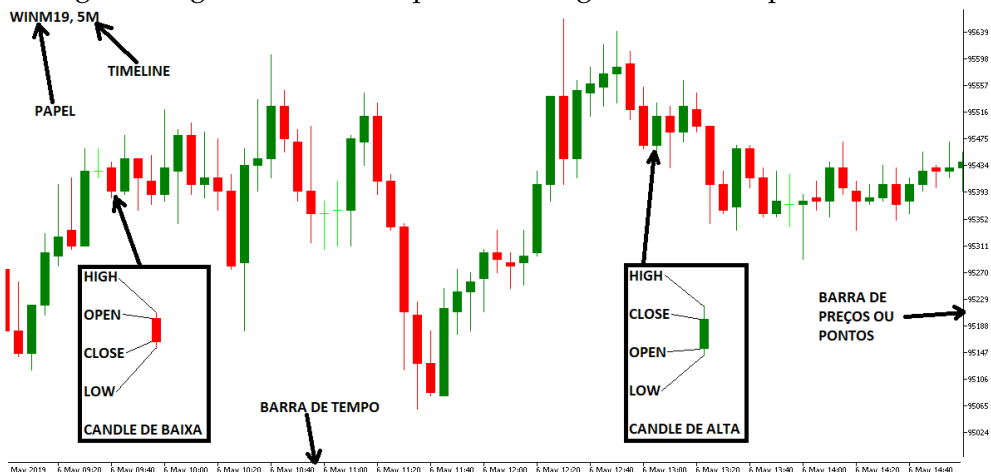
A métrica utilizada é o **Coefficiente de Correlação de 'Person'**, aplicado entre os conjuntos de registros/observações do gráfico de referência e de cada parte, de igual tamanho da referência, captada no gráfico de histórico no processo de varredura. A correlação de 'Person' é considerada o tipo normal dentre outras existentes e será utilizada por ser mais simples. E o valor da métrica, que servirá de fronteira entre as partes semelhantes ou não em relação ao gráfico de referência, será determinada pelo próprio usuário, lembrando que esse coeficiente é um número real entre '-1.0' e '1.0', significando que quanto mais próximo de '1.0', mais semelhante será.

Ao aplicar os cálculos de correlação, este projeto considera a comparação entre os pares de pontos na **forma bruta (raw)**, sem nenhuma transformação. Utilizaremos a linguagem de programação 'Python', com uso do ambiente de programação 'Jupyter Notebook'.

## 1.2 1. Conceitos

### 1.2.1 1.1 'Candlestick'

O gráfico de movimentações financeiras de uma ação de empresa ou índice de bolsa de valores pode ser do tipo 'candlestick' que tem como ponto básico o 'candle' por ser parecido com uma vela. A figura a seguir mostra uma parte de um gráfico desse tipo com a estrutura de 'candles'.



Onde: - '**HIGH**' é o valor máximo que o preço ou pontuação chegou dentro do tempo do 'candle'. - '**LOW**' é o valor mínimo que o preço ou pontuação chegou dentro do tempo do 'candle'. - '**OPEN**' é o valor de abertura de preço ou pontuação que iniciou o 'candle' atual ao terminar o tempo do 'candle' anterior. - '**CLOSE**' é o valor de fechamento de preço ou pontuação que terminou o 'candle' atual, dando início ao 'candle' seguinte. - '**CANDLE DE ALTA**' é o 'candle' onde o preço ou pontuação de fechamento ('**CLOSE**') é **MAIOR** que o de abertura ('**OPEN**'), informando que no tempo decorrido houve um aumento de valor. Neste caso, a alta no valor é representada pela cor verde do 'candle', mas não é padrão, podendo ser usada qualquer cor. - '**CANDLE DE BAIXA**' é o 'candle' onde o preço ou pontuação de fechamento ('**CLOSE**') é **MENOR** que o de abertura ('**OPEN**'), informando que no tempo decorrido houve uma diminuição de valor. Neste caso, a baixa no valor é representada pela cor vermelha do 'candle', mas não é padrão, podendo ser usada qualquer cor. - **PAPEL** é a identificação do produto financeiro, ou seja, o que está sendo negociado, podendo ser ações de empresas e índices da bolsa de valores. Neste caso, o papel '**WINM19**' significa que o gráfico é de negociação do Mini-Índice Futuro com vencimento em Junho de 2019, formado pela média das principais ações da **B3** (bolsa de valores brasileira, sediada em São Paulo - SP). Como outros exemplos, a B3 tem ações da Companhia Vale (VALE3), da Petrobrás (PETR3, PETR4), entre muitos outros. Maiores informações podem ser encontradas no portal da [B3](#). - '**TIMELINE**' é o padrão de tempo para todos os 'candle'. Neste caso, **5M** significa que cada 'candle' representa a variação de preço ou pontuação, da abertura ao fechamento, em exatamente 5 minutos. algumas das 'timelines' usadas, são 1min, 5min, 15min, 1h, 4h, 1d, 1sem e 1mes. - **BARRA DE TEMPO** é a escala temporal no eixo horizontal. Neste caso, cada 'candle' dista em 5 minutos de um para o outro na sequência temporal. - **BARRA DE PREÇOS OU PONTOS** é a escala, no eixo vertical, de variação de preços de ações de empresas ou de pontos de índices de bolsas de valores. Neste caso, são mostrados pontuações de índices.

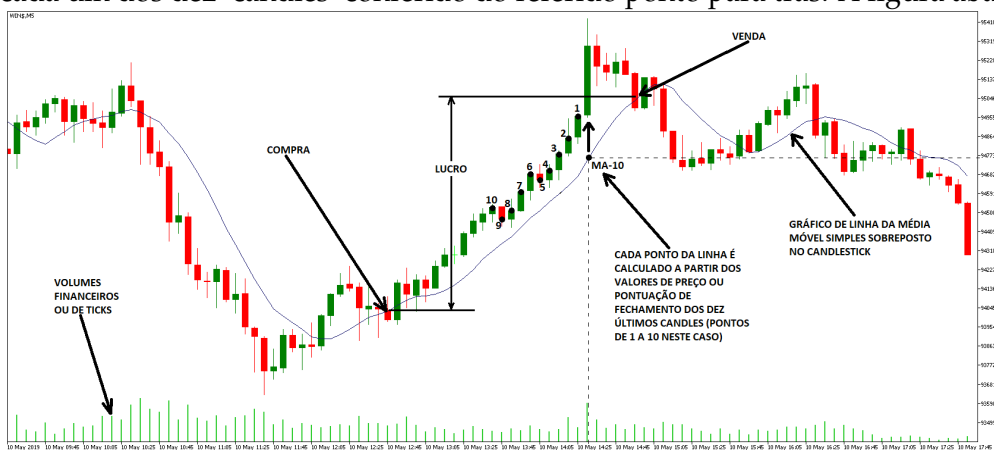
### 1.2.2 1.2 Plataforma de Negociação

É um software onde os operadores de bolsas de valores podem visualizar o gráfico e, mediante uma conta credenciada numa corretora que presta serviço a uma bolsa de valores, podem executar ordens de negociação de compra e venda de papéis. Há muitas plataformas e todas elas tem características comuns e próprias. Dentre as comuns, podemos alterar a 'timeline', as cores e o tipo de gráfico (há outros tipos além do candlestick, como de linha e OHLC).

Uma das plataformas de negociação mais difundidas no mundo é a **MetaTrader 5**, mantido pela empresa **MetaQuotes Software Corp.**. Seu instalador com instruções está disponível no site da [MetaQuotes](#) para o sistema operacional Windows. Usaremos o MetaTrader 5 para recuperarmos dados de um determinado papel usando uma conta do tipo 'demo' adquirido após cadastro em uma corretora brasileira que disponibiliza operações na B3 através dessa e de outras plataformas. Mas a própria MetaQuotes disponibiliza gratuitamente, sem a necessidade de uma corretora, uma conta 'demo' para qualquer usuário operar de forma simulada no mercado **FOREX** (tipo de mercado de abrangência mundial, semelhante a uma bolsa de valores, que disponibiliza operações de câmbio entre as principais moedas do mundo).

### 1.2.3 1.3 Indicadores e Robôs Investidores

Os **Indicadores** são ferramentas disponibilizadas pelas plataformas de negociação para auxiliar o operador na sua tomada de decisão. Por exemplo, temos a clássica **Média Móvel - MA (Movie Average)** que é o resultado da média aritmética dos 'N' últimos 'candles' considerando um dos preços padrões dos 'candles' (OPEN, HIGH, LOW, CLOSE), mas podendo ser usado outros critérios de cálculos para a média. Num exemplo hipotético, se considerarmos os dez últimos 'candles' e o preço 'CLOSE', será desenhado um gráfico de linha sobreposto ao 'candlestick' onde cada ponto da linha foi calculado pela média aritmética de todos os preços ou pontuações 'CLOSE' de cada um dos dez 'candles' conferido do referido ponto para trás. A figura abaixo mostra isso.



A média móvel, com suas variações, é só um dos inúmeros indicadores existentes. E a partir deles, podemos desenvolver os tão chamados **Robôs Investidores**. Na plataforma 'MetaTrader 5', eles são chamados de 'Expert Advisor - EA' e o usuário pode criar um EA baseado totalmente em um ou mais indicadores. A diferença entre indicador e robô é que a primeira é utilizada pelo operador para que ele mesmo dê entrada e saída de uma negociação. Por exemplo, numa tendência de subida do gráfico de 'candlestick', utilizando os devidos botões de comando da plataforma, o operador realiza uma ordem de compra em algum momento quando o preço ou pontuação atual toca na linha de média móvel e, posteriormente, realiza a venda no toque seguinte, concretizando algum lucro. A figura anterior também mostra essa negociação. Já o segundo, o robô, faz esse

processo de forma automática, sem intervenção do operador. É claro que desenvolver um robô não se resume a utilizar apenas indicadores, mas fica a critério do programador utilizar também código proprietário para otimizar as operações.

Devemos ter em mente no exemplo prático anterior, para melhor entendimento, que no momento do primeiro toque do preço ou pontuação atual com a média móvel, os 'candles' seguintes ainda serão formados e o operador deve prever para qual direção o valor irá, se cair ou subir, obtendo lucro ou prejuízo. E para otimizar isso, ele pode inclusive utilizar ferramentas estatísticas e até mesmo sofisticados algoritmos de inteligência artificial. Neste trabalho, utilizaremos apenas ferramentas estatísticas.

#### 1.2.4 1.4 Volumes de Negócios

Dentre os tipos de indicadores quanto à sua origem, há os calculados e os originais. Os próprios 'candles' podem ser vistos como indicadores originais, enquanto a média móvel é calculada. Outros indicadores originais, coletados diretamente dos servidores de dados, são os volumes de negócios, sendo dois tipos, e o 'spread'. Falaremos aqui dos volumes de negócios.

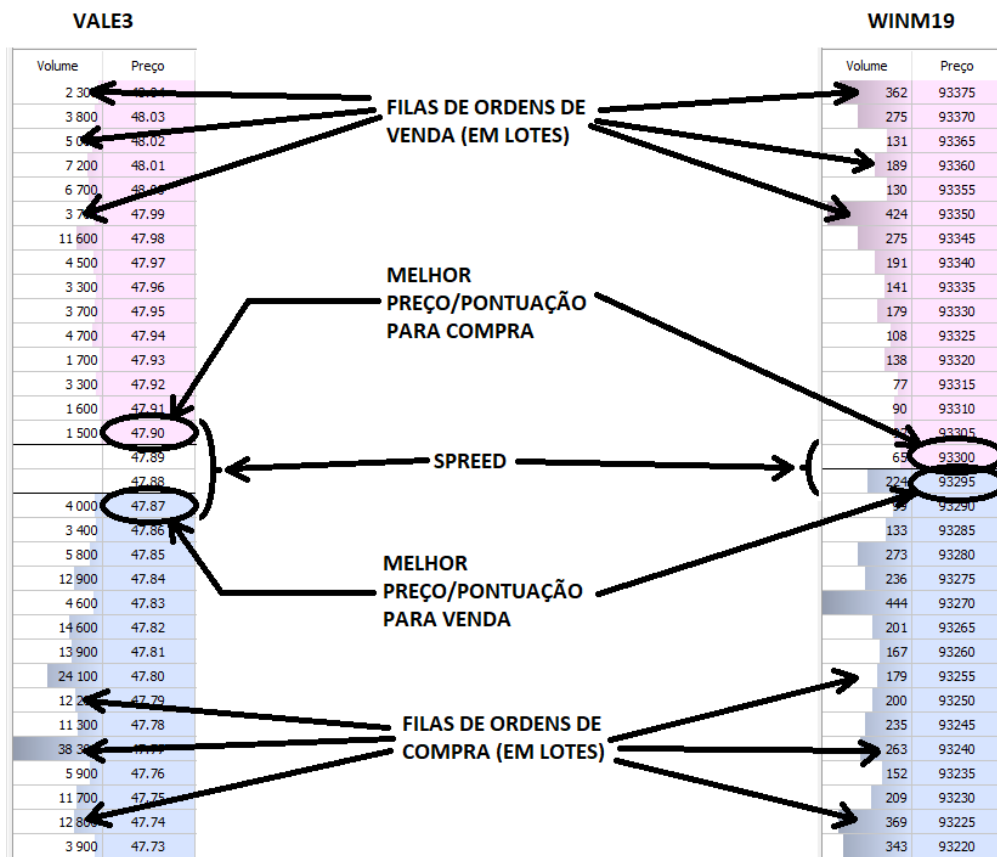
O primeiro tipo é o **Volume Financeiro**, no qual é registrado para cada 'candle' a soma de todos os **lotes** comprados e vendidos. O lote é a unidade básica de volume financeiro, sendo que cada papel tem seu valor. Por exemplo, o papel 'PETR4' da B3 tem o valor do lote a R\$26,62 e o valor para o papel 'WINM19' é de R\$25,00. Porém, no 'PETR4' só pode ser negociado múltiplos de 100 lotes, enquanto no 'WINM19' pode ser negociado qualquer quantidade. Mas essas regras de negociação não terão efeito neste projeto. O operador pode comprar e/ou vender a quantidade de lotes que quiser, dentro de seus limites financeiros, seguindo as regras de cada papel.

O segundo tipo é o **Volume de 'Ticks'**, no qual é registrado para cada 'candle' a soma de todos os 'ticks' comprados e vendidos. O 'tick' é a unidade básica de movimentação de preço ou pontuação. Para entender melhor, saiba que o preço ou pontuação se movimenta de forma discreta, ou seja, sobe ou desce 'tick' a 'tick'. Por exemplo, o preço no papel 'PETR4' varia a cada R\$0,01. Se o preço desse papel estiver num determinado momento a R\$21,43 e um operador executa uma ordem de compra imediata, o sistema da B3, obviamente, irá executar essa compra com quem está vendendo. Assim, o sistema irá procurar uma ordem de venda de outro operador que esteja na fila de um 'tick' mais próximo, ou seja, a R\$21,44, ou a R\$21,45, ou a R\$21,46, etc, sempre a cada R\$0,01. Para o papel 'WINM19', que utiliza pontuação em vez de preço, o 'tick' é de 5 pontos, ou seja, o valor do índice sobe ou desce a cada 5 pontos. Portanto, somando-se todos os 'ticks' movimentados, para cima e para baixo, dentro do tempo de um 'candle', temos o seu respectivo volume de 'ticks'. A figura anterior, do item '1.3', também mostra o indicador de volume, podendo ser financeiro ou de 'ticks'.

Todas essas informações podem ser encontradas com maiores detalhes na documentação da plataforma 'MetaTrader 5' e em muitos 'sites' especializados no assunto.

#### 1.2.5 1.5 Spread e Profundidade de Mercado

Todo papel negociado tem dois preços ou pontuações correntes: (melhor) compra e (melhor) venda. O primeiro será sempre maior que o segundo e a diferença entre eles nos dá o '**spread**', como mostra a figura abaixo. Nessa figura, temos uma tabela dinâmica extraída da plataforma 'MetaTrader 5' que é chamada de profundidade de mercado, referentes aos papéis 'VALE3' e 'WINM19'. As figuras representam o dinamismo das negociações paralizadas em um dado momento. Para entender melhor, saiba que o sistema tem, basicamente, dois tipos de ordens para compra e venda: direta (a mercado) e posicionada (na fila).



Quando um operador envia uma ordem de **compra direta**, também conhecida como ordem a mercado, com uma determinada quantidade de lotes, o sistema compra essa quantidade diretamente no melhor preço ou pontuação de compra, que na verdade é uma fila de lotes de venda, pois sempre compramos de quem está vendendo. Caso a quantidade de lotes na fila do melhor preço/pontuação de compra seja maior que a determinada pelo operador, apenas há a subtração e a quantidade na fila é reduzida. Mas, quando tem quantidade menor, a subtração é feita e a fila nesse preço/pontuação é zerada. Para completar a compra dos lotes restantes da ordem, o sistema continua realizando a compra no preço/pontuação seguinte de compra (acima), desde que tenha lotes de venda suficientes na fila, passando a ser o próximo melhor preço/pontuação para compra. Se ainda não tiver lotes suficientes, o sistema passa para o próximo preço/pontuação, repetindo esse processo até acabar com os lotes de compra da ordem do operador. O que descrevemos aqui foi o processo de movimento do preço/pontuação para cima (alta). O mesmo acontece de forma análoga e inversa para uma ordem direta de venda, ocorrendo o processo de movimento do preço/pontuação para baixo (baixa).

Quando o operador envia uma ordem de **compra posicionada**, não há de fato uma execução da ordem, mas um posicionamento dos lotes em fila num preço ou pontuação determinado pelo próprio operador. Por exemplo, utilizando as informações da figura anterior, se ele enviar uma ordem de compra com 100 lotes para o papel 'VALE3' no preço R\\$/47,76, esses lotes ficarão na fila do preço mencionado e se somarão aos 5900 lotes já existentes lá, resultando em 6000 lotes. O posicionamento significa que o operador espera que o melhor preço de compra do papel caia do atual R\\$/47,90 para R\\$/47,76, quando o sistema executará, de fato, a compra dos 100 lotes que ele quer. Diferentemente da ordem direta, na posicionada não há movimento de preço caso a quantidade de lotes de compra na fila seja maior que a quantidade consumida (execução de

compra de quem está vendendo). Continuando no exemplo, se o melhor preço caiu até chegar na posição desejada do operador e, momentaneamente, houver uma execução que consumiu 70 lotes da ordem com o melhor preço subindo logo em seguida, a ordem ainda terá 30 lotes de compra na fila do preço desejado. Nesse caso, o operador deverá esperar o preço baixar novamente até que todos os lotes dele sejam consumidos. Todo esse processo acontece de forma análoga e inversa para uma ordem posicionada de venda.

Sobre o 'spread', é uma característica muito ligada ao dinamismo do preço/pontuação e liquidez (volume financeiro) do papel e é medida em 'ticks'. Há papéis que raramente o 'spread' é alto, como no 'WINM19' da figura anterior ( $93300 - 93295 = 5 = 1$  'tick'). Já no 'VALE3', tem momentos de menor dinamismo, como também mostrado na figura anterior ( $47,90 - 47,87 = 0,03 = 3$  'ticks'). Há papéis com 'spreads' muito altos, com mais de 5 'ticks', geralmente com liquidez muito baixa, significando pouco volume de negociação.

Todas essas informações também podem ser encontradas com maiores detalhes na documentação da plataforma 'MetaTrader 5' e em muitos 'sites' especializados no assunto.

### 1.3 2. Base de Dados

A base de dados ('dataset') utilizada neste projeto foi extraída por uma ferramenta de recuperação de dados históricos da plataforma 'MetaTrader 5' que faz um download a partir da base de dados da B3 por intermédio dos servidores de uma corretora. Esses dados são específicos do papel 'WINM19', com 'timeline' de **5min**, com dezenas de milhares de 'candles' com seus respectivos preços de abertura, máximo, mínimo e de fechamento coletados entre as datas-horas '02-04-2019 09:00:00' e '30-04-2019 11:35:00', sempre dentro do horário de expediente da B3 para esse papel que ocorre nos dias úteis entre 09:00:00 e 18:00:00. Vejamos um 'overview' dos dados:

```
[1]: # Carregando os dados do arquivo para um 'dataframe'...

import pandas as pd

df_raw = pd.read_csv('B3_CANDLES_WIN_M5_201804020900_201904301135.csv',
                    sep='\t')

# Mostrando os cinco primeiros registros do 'dataframe'...
print('Dimensões:', df_raw.shape)
df_raw.head()
```

Dimensões: (28580, 9)

```
[1]:
```

	<DATE>	<TIME>	<OPEN>	<HIGH>	<LOW>	<CLOSE>	<TICKVOL>	<VOL>	\
0	2018.04.02	09:00:00	91081	91225	90942	90995	5494	20340	
1	2018.04.02	09:05:00	90995	91017	90937	90979	4640	15613	
2	2018.04.02	09:10:00	90979	91097	90969	91070	4679	13901	
3	2018.04.02	09:15:00	91065	91102	91038	91043	3605	11205	
4	2018.04.02	09:20:00	91043	91065	91001	91043	3565	10508	

```
<SPREAD>
```

0	1
1	1

```
2      1
3      1
4      1
```

[2]: *# Mostrando os cinco últimos registros do 'dataframe'...*

```
df_raw.tail()
```

[2]:

	<DATE>	<TIME>	<OPEN>	<HIGH>	<LOW>	<CLOSE>	<TICKVOL>	<VOL>	\
28575	2019.04.30	11:15:00	96900	96935	96715	96725	20861	70833	
28576	2019.04.30	11:20:00	96730	96850	96640	96670	27967	97054	
28577	2019.04.30	11:25:00	96665	96750	96500	96560	27718	95275	
28578	2019.04.30	11:30:00	96565	96660	96525	96645	12092	41762	
28579	2019.04.30	11:35:00	96645	96680	96425	96525	25128	89698	

	<SPREAD>
28575	1
28576	1
28577	1
28578	1
28579	1

O 'dataframe' extraído diretamente do arquivo apresenta as seguintes 'features': - < DATE > : Data de criação do 'candle'. - < TIME > : Hora de criação do 'candle'. - < OPEN > : Valor (pontuação) de abertura do 'candle'. - < HIGH > : Valor (pontuação) máximo do 'candle'. - < LOW > : Valor (pontuação) mínimo do 'candle'. - < CLOSE > : Valor (pontuação) de fechamento do 'candle'. - < TICKVOL > : Volume de ticks negociado no período do 'candle'. - < VOL > : Volume financeiro negociado no período do 'candle'. - < SPREAD > : Diferença entre os melhores valores (pontuação) de compra e venda no fechamento do 'candle'.

Cada observação (linha de registro) corresponde a um 'candle' de 5min. Vamos agora reorganizar os nomes das 'features' removendo os sinais 'maior que' e 'menor que', bem como deixar mais adequado o nome 'VOL' alterando para 'FINVOL' (volume financeiro).

[3]: *# Alterando os nomes das 'features'...*

```
df_raw.columns = ['DATE', 'TIME', 'OPEN', 'HIGH', 'LOW', 'CLOSE', 'TICKVOL', 'FINVOL', 'SPREAD']
df_raw.head()
```

[3]:

	DATE	TIME	OPEN	HIGH	LOW	CLOSE	TICKVOL	FINVOL	SPREAD
0	2018.04.02	09:00:00	91081	91225	90942	90995	5494	20340	1
1	2018.04.02	09:05:00	90995	91017	90937	90979	4640	15613	1
2	2018.04.02	09:10:00	90979	91097	90969	91070	4679	13901	1
3	2018.04.02	09:15:00	91065	91102	91038	91043	3605	11205	1
4	2018.04.02	09:20:00	91043	91065	91001	91043	3565	10508	1

Vamos observar se os dados estão devidamente tipificados.

[4]: *# Mostra os tipos de dados de todas as 'features'...*

```
print(df_raw.dtypes)
```

```
DATE          object
TIME          object
OPEN          int64
HIGH          int64
LOW           int64
CLOSE         int64
TICKVOL       int64
FINVOL        int64
SPREAD        int64
dtype: object
```

Percebemos que 'DATE' e 'TIME' não estão de fato como de tipos data e hora. Portanto, vamos alterá-los e reorganizar o 'dataframe' para que as duas 'features' sejam uma só e do tipo 'data-hora'.

[5]: *# Alterando tipos de dados e 'features' novamente...*

```
df_raw.DATE = list(str(a) + ' ' + str(b) for a, b in zip(df_raw.DATE.values,
→df_raw.TIME.values))
df_raw = df_raw.drop(['TIME'], axis=1)
df_raw.rename(columns={'DATE': 'DATETIME'}, inplace=True)
df_raw.DATETIME = pd.to_datetime(df_raw.DATETIME)
print(df_raw.dtypes)
df_raw.head()
```

```
DATETIME      datetime64[ns]
OPEN          int64
HIGH          int64
LOW           int64
CLOSE         int64
TICKVOL       int64
FINVOL        int64
SPREAD        int64
dtype: object
```

[5]:

	DATETIME	OPEN	HIGH	LOW	CLOSE	TICKVOL	FINVOL	SPREAD
0	2018-04-02 09:00:00	91081	91225	90942	90995	5494	20340	1
1	2018-04-02 09:05:00	90995	91017	90937	90979	4640	15613	1
2	2018-04-02 09:10:00	90979	91097	90969	91070	4679	13901	1
3	2018-04-02 09:15:00	91065	91102	91038	91043	3605	11205	1
4	2018-04-02 09:20:00	91043	91065	91001	91043	3565	10508	1

---

Agora vamos observar os dados estatísticos básicos.

[6]: *# Mostrando números estatísticos gerais do 'dataframe'...*

```
df_raw.describe()
```

[6]:

	OPEN	HIGH	LOW	CLOSE \
count	28580.000000	28580.000000	28580.000000	28580.000000
mean	88595.329986	88677.080651	88512.41627	88594.803079
std	7300.254217	7297.881916	7301.53694	7300.122394
min	72443.000000	72752.000000	72417.000000	72741.000000
25%	82000.000000	82076.000000	81921.000000	82000.000000
50%	89445.000000	89534.500000	89358.000000	89442.000000
75%	95320.000000	95391.000000	95250.000000	95320.000000
max	101802.000000	101827.000000	101732.000000	101802.000000

	TICKVOL	FINVOL	SPREAD
count	28580.000000	28580.000000	28580.0
mean	10562.033345	36113.442022	1.0
std	6749.222027	22402.810567	0.0
min	375.000000	1631.000000	1.0
25%	5825.000000	20326.750000	1.0
50%	9096.000000	31358.000000	1.0
75%	13649.250000	46610.000000	1.0
max	74144.000000	248721.000000	1.0

[7]: *# Verificando dados faltantes...*

```
df_raw.isnull().sum()
```

[7]:

DATETIME	0
OPEN	0
HIGH	0
LOW	0
CLOSE	0
TICKVOL	0
FINVOL	0
SPREAD	0

dtype: int64

---

Percebemos que não há dados faltantes e que a 'feature' 'SPREAD' tem todas as amostras iguais a '1.0'. Portanto, podemos eliminá-la do 'dataframe'.

[8]: *# Removendo a 'feature' 'SPREAD'...*

```
df_raw = df_raw.drop(['SPREAD'], axis=1)
df_raw.head()
```

```
[8]:
```

		DATETIME	OPEN	HIGH	LOW	CLOSE	TICKVOL	FINVOL
0	2018-04-02	09:00:00	91081	91225	90942	90995	5494	20340
1	2018-04-02	09:05:00	90995	91017	90937	90979	4640	15613
2	2018-04-02	09:10:00	90979	91097	90969	91070	4679	13901
3	2018-04-02	09:15:00	91065	91102	91038	91043	3605	11205
4	2018-04-02	09:20:00	91043	91065	91001	91043	3565	10508

## 1.4 Aplicando a Correlação de Person nos Dados Brutos

O que queremos ao comparar dois conjuntos de pontos, formando gráficos, é verificar se ambos tem variação semelhante, independente da distância entre os pares de valores de cada ponto dos conjuntos. Vamos usar a função 'pearsonr()' do módulo 'stats' da biblioteca 'scipy' que retorna uma tupla com o coeficiente de correlação e a probabilidade de ser ao acaso, nessa ordem, no qual iremos utilizar somente o primeiro.

A estrutura do processo de varredura no 'dataframe' para cálculo da correlação, requer o seguinte: - Um **'dataset' de dados brutos**, representando todo o **histórico do 'candlestick'**, por onde nosso sistema fará a varredura 'candle' a 'candle' na sequência temporal. - Um **'dataset' de referência**, bem menor, representando a(s) parte(s) que queremos encontrar no dataset de dados brutos. Chamaremos de **'subset' padrão**. - **'Datasets' de comparação**, sendo todos do mesmo tamanho do 'subset' padrão, extraídos 'candle' a 'candle' na varredura do 'dataset' de dados brutos. Chamaremos de **'subsets' extraídos**. - Uma **referência de pontuação** de cada 'candle' para comporem os elementos dos 'subsets'. Podemos utilizar qualquer uma das 'features', mas neste projeto vamos nos concentrar na **pontuação de fechamento - 'CLOSE'**, pois é a mais utilizada pelos operadores de bolsas de valores.

Para começar, vamos escolher, hipoteticamente, um 'subset' padrão. Poderíamos extrair de outro arquivo baixado da B3 pela plataforma 'MetaTrader 5', fora do intervalo de datas do arquivo que gerou nosso 'dataframe', mas vamos extrair do próprio 'dataset' de dados brutos. Isso nos dá a convicção de que um dos 'subsets' extraídos na varredura deve ser exatamente o mesmo do padrão, obtendo coeficiente de correlação igual a '1.0'. Os códigos a seguir mostram a parte do gráfico de 'candlestick' onde extraímos nosso 'subset' padrão, usando a biblioteca 'plotly' que incorpora ferramentas para gráficos de 'candlestick'. Vamos analisá-lo.

```
[9]: # Bibliotecas necessárias para o tratamento de gráficos do tipo 'candlestick'...
```

```
import plotly.offline as py
import plotly.graph_objs as go
import matplotlib.pyplot as plt
from datetime import datetime
import numpy as np
py.init_notebook_mode()
```

```
[10]: # Extração de uma parte do grande 'dataset' de dados brutos...
```

```
df_raw_part = df_raw[np.logical_and(df_raw['DATETIME'] >= '2018-10-23 09:45:00',
                                     df_raw['DATETIME'] <= '2018-10-25 15:55:
→00')]
print('Dimensões:', df_raw_part.shape)
df_raw_part.head()
```

Dimensões: (289, 7)

```
[10]:
```

	DATETIME	OPEN	HIGH	LOW	CLOSE	TICKVOL	FINVOL
15096	2018-10-23 09:45:00	87608	87680	87536	87572	10401	28198
15097	2018-10-23 09:50:00	87572	87587	87443	87459	9001	23003
15098	2018-10-23 09:55:00	87459	87505	87289	87351	14013	36648
15099	2018-10-23 10:00:00	87351	87361	87243	87269	12587	36902
15100	2018-10-23 10:05:00	87269	87366	87259	87351	11391	33163

```
[11]: # Gráfico de 'candlestick' tomando como base temporal a 'feature' 'DATETIME'...
```

```
trace = go.Candlestick(x=df_raw_part['DATETIME'],
                       open=df_raw_part['OPEN'],
                       high=df_raw_part['HIGH'],
                       low=df_raw_part['LOW'],
                       close=df_raw_part['CLOSE'])

layout = go.Layout(xaxis = dict(rangeslider = dict(visible = False)))

data = [trace]

fig = go.Figure(data=data,layout=layout)
py.offline.iplot(fig, filename='simple_candlestick_without_range_slider')
```



Tomando como base temporal, no eixo horizontal, a 'feature' 'DATETIME', observamos claramente a falta de registros de 'candles' fora do horário de expediente de trabalho da bolsa de valores. Isso não é bom para nosso trabalho. Portanto, iremos considerar, a partir de agora, o índice do 'dataset' de dados brutos como nossa base temporal, eliminando os espaços vazios fora do

expediente, mas mantendo o mesmo intervalo de data e hora. Então, podemos escolher nosso 'subset' padrão.

```
[12]: # Gráfico de 'candlestick' tomando como base temporal o índice do 'dataset' de
      ↪ dados brutos
      # e mostrando o nosso 'subset' padrão...

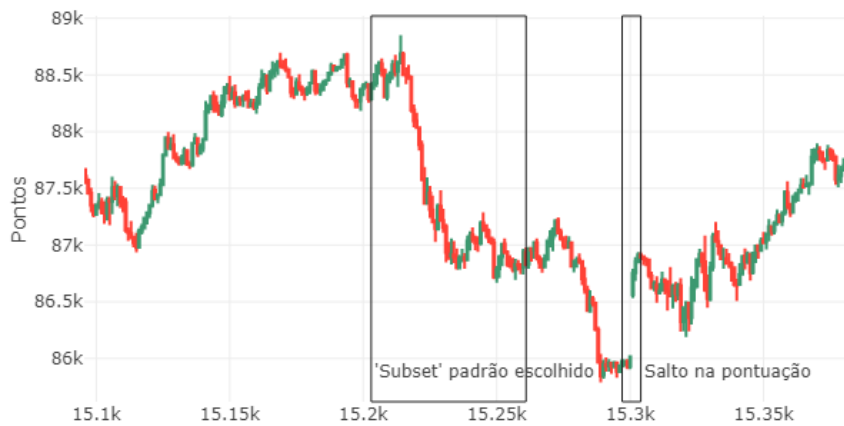
trace = go.Candlestick(x=df_raw_part.index,
                       open=df_raw_part['OPEN'],
                       high=df_raw_part['HIGH'],
                       low=df_raw_part['LOW'],
                       close=df_raw_part['CLOSE'])

layout = {'title': "Parte do 'dataset' e escolha do 'subset' padrão",
          'yaxis': {'title': 'Pontos'},
          'xaxis': {'rangeslider': {'visible': False}},
          'shapes': [{ 'x0': 15297,
                        'x1': 15304,
                        'y0': 0, 'y1': 1, 'xref': 'x', 'yref': 'paper',
                        'line': {'color': 'rgb(30,30,30)', 'width': 1}},
                     { 'x0': 15203,
                        'x1': 15261,
                        'y0': 0, 'y1': 1, 'xref': 'x', 'yref': 'paper',
                        'line': {'color': 'rgb(30,30,30)', 'width': 1}}
                    ],
          'annotations': [{ 'x': 15304,
                              'y': 0.05,
                              'xref': 'x', 'yref': 'paper',
                              'showarrow': False, 'xanchor': 'left',
                              'text': 'Salto na pontuação'},
                           { 'x': 15203,
                              'y': 0.05,
                              'xref': 'x', 'yref': 'paper',
                              'showarrow': False, 'xanchor': 'left',
                              'text': "'Subset' padrão escolhido"}
                          ]
        }

data = [trace]

fig = go.Figure(data=data,layout=layout)
py.offline.iplot(fig, filename='simple_candlestick_without_range_slider')
```

Parte do 'dataset' e escolha do 'subset' padrão



Nessa parte do 'dataset', eliminamos os espaços vazios em relação ao eixo horizontal, mas no eixo vertical percebemos um salto de pontuação do índice '15300' para o '15301'. Isso acontece porque nas primeiras negociações no início do expediente de trabalho o preço ou pontuação negociado está distante do último valor do dia anterior. Nesse papel, esses saltos são comuns e vamos aceitar normalmente em nosso projeto. Abaixo, temos nosso 'subset' padrão que utilizaremos em todo nosso trabalho.

```
[13]: # Nosso 'subset' padrão...

subset_default = df_raw_part.loc[15203:15261]
print('Dimensões:', subset_default.shape)
subset_default.head()
```

Dimensões: (59, 7)

```
[13]:
```

	DATETIME	OPEN	HIGH	LOW	CLOSE	TICKVOL	FINVOL
15203	2018-10-24 09:45:00	88291	88435	88271	88415	10671	30563
15204	2018-10-24 09:50:00	88415	88517	88368	88430	14593	39926
15205	2018-10-24 09:55:00	88430	88589	88394	88584	10034	29821
15206	2018-10-24 10:00:00	88579	88625	88502	88610	11614	34996
15207	2018-10-24 10:05:00	88599	88646	88486	88543	11067	34382

```
[14]: # Gráfico de 'candlestick' do nosso 'subset' padrão...

trace = go.Candlestick(x=subset_default.index,
                        open=subset_default['OPEN'],
                        high=subset_default['HIGH'],
                        low=subset_default['LOW'],
```

```

close=subset_default['CLOSE'])

layout = go.Layout(title = "Amostra ampliada do 'Subset' padrão escolhido",
                    xaxis = dict(rangeslider = dict(visible = False)))

data = [trace]

fig = go.Figure(data=data,layout=layout)
py.offline.iplot(fig, filename='simple_candlestick_without_range_slider')

```



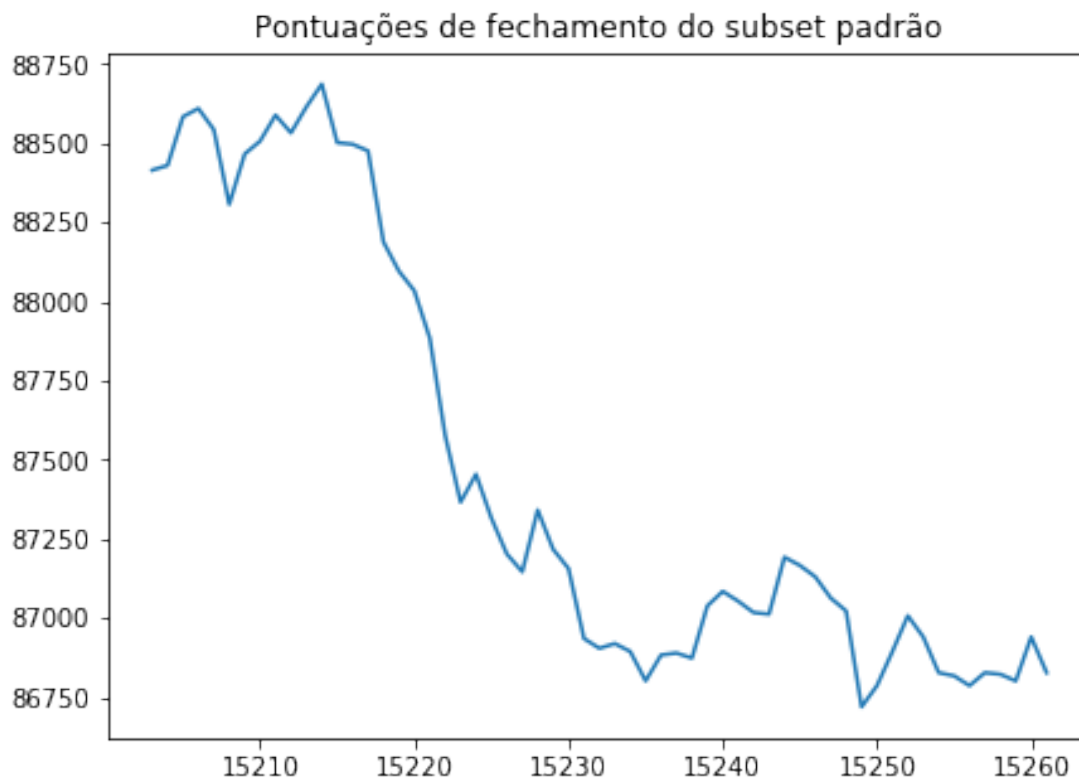
Agora, vamos trabalhar na aplicação da correlação. Como já foi dito, escolhemos trabalhar com a pontuação de fechamento ('CLOSE'). Assim, podemos plotar um gráfico de linha somente com essas pontuações do 'subset' padrão.

[15]: *# Gráfico de linha das pontuações de fechamento do 'subset' padrão...*

```

fig = plt.figure(figsize=(7,5))
plt.plot(subset_default.index, subset_default.CLOSE)
plt.title("Pontuações de fechamento do subset padrão")
plt.show()

```



Determinamos, então, uma função que faz a varredura em todo o 'dataset', comparando cada 'subset' extraído com o padrão, retornando aqueles com coeficiente de correlação acima de um determinado valor. Para esse processo, seguiremos as seguintes premissas: - Criação de um parâmetro de fronteira para o coeficiente de correlação com o objetivo de separar os 'subsets' extraídos que sejam semelhantes dos não semelhantes ao padrão, caso seus coeficientes estejam acima ou abaixo do valor desse parâmetro. - Para o cálculo da correlação, evitando dados em branco, os 'subsets' extraídos devem ter exatamente o mesmo tamanho do padrão. - Vamos considerar que a função irá receber a 'feature' que se quer referenciar nos cálculos, não obrigatoriamente sendo a 'CLOSE', dando mais flexibilidade. - A varredura será feita do primeiro ao último 'candle' possível do 'dataset'. Uma ilustração de parte desse passo é mostrada no gráfico abaixo. Devemos considerar que o último 'candle' da varredura é aquele que coincide com a posição do último menos o tamanho do 'subset' padrão para que não tenhamos 'subsets' extraídos com tamanho menor. - A função retornará uma tupla com dois elementos. O primeiro será uma lista contendo as posições iniciais no 'dataset' onde foram encontradas 'subsets' extraídos semelhantes ao padrão. E o segundo elemento será uma lista dos coeficientes de correlação de 'Person' correspondente, na mesma ordem, a cada posição da lista do primeiro elemento.

[16]: *# Demonstração gráfica da varredura de 'subsets' extraídos...*

```
df_raw_part2 = df_raw_part.loc[15096:15170]
trace = go.Candlestick(x=df_raw_part2.index,
```

```

        open=df_raw_part2['OPEN'],
        high=df_raw_part2['HIGH'],
        low=df_raw_part2['LOW'],
        close=df_raw_part2['CLOSE'])

layout = {'title': "Demonstração gráfica da varredura de 'subsets' extraídos",
          'yaxis': {'title': 'Pontos'},
          'xaxis': {'rangeslider': {'visible': False}},
          'shapes': [{ 'x0': 15119, 'x1': 15146, 'y0': 0.05, 'y1': 0.9, 'xref': 'x', 'yref': 'paper',
                        'line': {'color': 'rgb(30,30,30)', 'width': 1}},
                      { 'x0': 15120, 'x1': 15147, 'y0': 0.1, 'y1': 0.95, 'xref': 'x', 'yref': 'paper',
                        'line': {'color': 'rgb(30,30,30)', 'width': 1}},
                      { 'x0': 15121, 'x1': 15148, 'y0': 0.15, 'y1': 1.0, 'xref': 'x', 'yref': 'paper',
                        'line': {'color': 'rgb(30,30,30)', 'width': 1}}
                    ],
          'annotations': [{ 'x': 15117, 'y': 0.5, 'xref': 'x', 'yref': 'paper',
                             'showarrow': False, 'xanchor': 'left',
                             'text': '...',
                             { 'x': 15119, 'y': 0.05, 'xref': 'x', 'yref': 'paper',
                               'showarrow': False, 'xanchor': 'left',
                               'text': 'subset n'},
                             { 'x': 15120, 'y': 0.10, 'xref': 'x', 'yref': 'paper',
                               'showarrow': False, 'xanchor': 'left',
                               'text': 'subset n+1'},
                             { 'x': 15121, 'y': 0.15, 'xref': 'x', 'yref': 'paper',
                               'showarrow': False, 'xanchor': 'left',
                               'text': 'subset n+2'},
                             { 'x': 15148, 'y': 0.5, 'xref': 'x', 'yref': 'paper',
                               'showarrow': False, 'xanchor': 'left',
                               'text': '...' }
                           ]
        }

data = [trace]

fig = go.Figure(data=data,layout=layout)
py.offline.iplot(fig, filename='simple_candlestick_without_range_slider')

```

### Demonstração gráfica da varredura de 'subsets' extraídos



[17]: *# Função de extração de 'subsets' semelhantes a um padrão.*

```
def get_similar_subsets(dataset, subset, col_dataset='CLOSE',  
    →col_subset='CLOSE', min_rate_corr=0.9, show_result=False):  
  
    """  
    Esta função faz uma busca por conjuntos menores de 'candles' em sequência,  
    →numa base maior  
    conforme um conjunto padrão. A métrica de correlação de 'Pearson' é  
    →utilizada como fronteira  
    para dividir os conjuntos de 'candles' em semelhante ou não ao padrão. O  
    →retorno de valor  
    desta função é uma tupla com duas listas, sendo a primeira representando as  
    →posições no  
    'dataset' dos conjuntos semelhantes ao padrão e a segunda representando os  
    →coeficientes de  
    correlação deles.  
  
    VARIÁVEIS DE ENTRADA:  
  
    - dataset: tipo 'pandas.DataFrame'  
      É um 'dataframe', representando um conjunto de 'candles', onde será  
    →extraída as posições  
      que contenham subconjuntos semelhantes ao padrão.  
  
    - subset: tipo 'pandas.DataFrame'  
      É um 'dataframe', representando um subconjunto de 'candles', que  
    →servirá de padrão de  
      semelhança para os subconjuntos extraídos do 'dataset'.
```

- *col\_dataset: tipos 'str' ou 'int'*  
*É o nome ou número da coluna do 'dataset' para ser usada no cálculo da*  
*→ correlação.*
  
- *col\_subset: tipos 'str' ou 'int'*  
*É o nome ou número da coluna do 'subset' para ser usada no cálculo da*  
*→ correlação.*
  
- *min\_rate\_corr: tipos 'float' ou 'int' no espaço [-1.0,1.0]*  
*É o parâmetro que representa o mínimo valor do coeficiente de*  
*→ correlação de 'Pearson' que*  
*servirá de fronteira para os 'subsets' extraídos, sendo que cada um*  
*→ terá seu coeficiente*  
*calculado, onde esse parâmetro irá classificá-los internamente em*  
*→ semelhante ou não ao*  
*padrão, separando os semelhantes como resultado da função.*
  
- *show\_result: tipo 'bool'*  
*É um ativador que serve para imprimir ou não na saída o resultado.*

#### *RESULTADO:*

- *return: tipo 'tuple(list(),list())'*  
*Esta função retorna uma tupla com dois elementos. O primeiro é uma*  
*→ lista contendo as*  
*posições (índices do 'dataset') que representam o primeiro 'candle' de*  
*→ cada 'subset'*  
*semelhante extraído. O segundo é uma lista, na mesma ordem da primeira,*  
*→ que representa os*  
*coeficientes de correlação de 'Pearson' de todos os 'subsets'*  
*→ semelhantes.*

"""

*## Tratamento de erros das variáveis de entrada.*

*import pandas*

*# Verificação de tipos.*

```
if type(dataset) != pandas.DataFrame:
    raise ValueError("'dataset' deve ser do tipo 'pandas.DataFrame'.")
if type(subset) != pandas.DataFrame:
    raise ValueError("'subset' deve ser do tipo 'pandas.DataFrame'.")
if not(type(col_dataset) == str or type(col_dataset) == int):
    raise ValueError("'col_dataset' deve ser do tipo 'str' ou 'int'.")
if not(type(col_subset) == str or type(col_subset) == int):
```

```

        raise ValueError("'col_subset' deve ser do tipo 'str' ou 'int'.")
    if not(type(min_rate_corr) == int or type(min_rate_corr) == float):
        raise ValueError("'rate_corr' deve ser do tipo 'int' ou 'float'.")

    # Verificação de conteúdo.
    if subset.shape[0] > dataset.shape[0]:
        raise ValueError("Tamanho do 'subset' (quantidade de registros) deve
→ser menor ou igual do 'dataset'.")
    if type(col_dataset) == str:
        try:
            float(dataset[col_dataset].values[0])
        except:
            raise ValueError("Os dados da coluna " + col_dataset + " do
→'dataset' devem ser numéricos.")
            if not(col_dataset in dataset.columns.values):
                raise ValueError("Não há nenhuma coluna no 'dataset' com o nome '"
→+ col_dataset + "'.")
            else:
                try:
                    float(dataset[dataset.columns[col_dataset]].values[0])
                except:
                    raise ValueError("Os dados da coluna " + col_dataset + " do
→'dataset' devem ser numéricos.")
                    if col_dataset < 0:
                        raise ValueError("'col_dataset' deve ser maior ou igual a zero.")
                    if col_dataset >= len(dataset.columns.values):
                        raise ValueError("Não há nenhuma coluna no 'dataset' com o índice
→'" + str(col_dataset) + "'.")
                    if type(col_subset) == str:
                        try:
                            float(subset[col_subset].values[0])
                        except:
                            raise ValueError("Os dados da coluna " + col_subset + " do 'subset'
→devem ser numéricos.")
                            if not(col_subset in subset.columns.values):
                                raise ValueError("Não há nenhuma coluna no 'subset' com o nome '"
→col_subset + "'.")
                            else:
                                try:
                                    float(subset[subset.columns[col_subset]].values[0])
                                except:
                                    raise ValueError("Os dados da coluna " + col_subset + " do 'subset'
→devem ser numéricos.")
                                    if col_subset < 0:
                                        raise ValueError("'col_subset' deve ser maior ou igual a zero.")
                                    if col_subset >= len(subset.columns.values):

```

```

        raise ValueError("Não há nenhuma coluna no 'subset' com o índice '" +
→+ str(col_subset) + "'.")
    if float(min_rate_corr) < -1.0 or float(min_rate_corr) > 1.0:
        raise ValueError("'rate_corr' deve estar no intervalo [-1.0,1.0].")

    ## Corpo principal da função.

    # Pegando os índices das colunas caso as variáveis de entrada 'col_dataset'
→e 'col_subset' sejam 'strings'.
    if type(col_dataset) == str:
        col_dataset = list(dataset.columns.values).index(col_dataset)
    if type(col_subset) == str:
        col_subset = list(dataset.columns.values).index(col_subset)

    # Carregamento das bibliotecas necessárias para cálculo de correlações.
    from scipy import stats
    from time import time
    from tqdm import tqdm

    # Varredura no 'dataset' para encontrar 'subsets' semelhantes ao padrão
→através do cálculo da correlação
    # de 'Pearson' entre as colunas dadas em 'col_dataset' e 'col_subset' de
→cada dataframe de entrada.
    size_subset = subset.shape[0]
    list_subset_similar = []
    list_corr_subset = []
    last_position = dataset.shape[0] - 1 - size_subset
    tm = time()
    for i in tqdm(range(last_position + 1)):
        if i > last_position:
            break
        subset_extracted = dataset.iloc[i:i+size_subset,:]
        corr,_ = stats.pearsonr(subset_extracted[subset_extracted.
→columns[col_dataset]].values,
                                subset[subset.columns[col_subset]].values)
        if corr >= min_rate_corr:
            list_subset_similar.append(dataset.index.values[i])
            list_corr_subset.append(corr)
    if show_result:
        print('- Varredura executada em', round(time()-tm,2), '\b(s).')
        print('- Foram encontrados', len(list_subset_similar), '\b(subsets)\b'
→extraídos semelhantes ao padrão',
                'utilizando as colunas ', ""+dataset.columns[col_dataset]+""",
→'do \b(dataset\b' e',
                ""+subset.columns[col_subset]+""", 'do \b(subset\b' com
→coeficientes de correlação de',

```

```
'\'Pearson\' igual ou maior que \'' + str(min_rate_corr) + '\'.')  
  
return list_subset_similar, list_corr_subset
```

---

A função no código acima é genérica, flexível. Não tem obrigatoriedade de que os dataframes sejam representações de 'candlesticks', podendo ser de qualquer natureza, desde que as colunas escolhidas obedeam as regras de tipos das variáveis de entrada. Neste trabalho, vamos utilizar essa função para usar a 'feature' 'CLOSE' do nosso 'dataset' de dados brutos e do 'subset' padrão que escolhemos. Como primeiro exemplo, vamos aplicar a função diretamente com os parâmetros padrões (colunas com 'CLOSE' e coeficiente de correlação mínima de '0.9'), mas mostrando o resultado na saída.

```
[18]: subsets, corrs = get_similar_subsets(df_raw, subset_default, show_result=True)
```

```
100%||  
28521/28521 [00:06<00:00, 4264.43it/s]
```

- Varredura executada em 6.7(s).
- Foram encontrados 627 'subsets' extraídos semelhantes ao padrão utilizando as colunas 'CLOSE' do 'dataset' e 'CLOSE' do 'subset' com coeficientes de correlação de 'Pearson' igual ou maior que '0.9'.

---

Nessa primeira tentativa, obtivemos uma quantidade muito grande de 'subsets' reconhecidos, dificultando a análise. Vamos repetir o processo aumentando o valor mínimo de coeficiente de correlação para reduzir essa quantidade.

```
[19]: subsets, corrs = get_similar_subsets(df_raw, subset_default, min_rate_corr=0.  
→95, show_result=True)
```

```
100%||  
28521/28521 [00:06<00:00, 4271.37it/s]
```

- Varredura executada em 6.68(s).
- Foram encontrados 30 'subsets' extraídos semelhantes ao padrão utilizando as colunas 'CLOSE' do 'dataset' e 'CLOSE' do 'subset' com coeficientes de correlação de 'Pearson' igual ou maior que '0.95'.

---

Agora, com bem menos 'subsets' reconhecidos, vamos mostrar se de fato são semelhantes ao padrão. Para simplificar, vamos utilizar o gráfico de linhas com apenas os dados da 'feature' 'CLOSE', lembrando que a lista do resultado da função anterior não são conjuntos inteiros de 'candles', mas apenas o índice no 'dataset' de dados brutos para o primeiro dos 'candles' de cada 'subset' reconhecido dentro desse 'dataset'.

[20]: *# Confirmando os dados do resultado...*

```
print("Índices em 'df_raw' de localização inicial dos 'subsets' reconhecidos:")
print(subsets)
print("-----")
print("Coeficientes de correlação, na mesma ordem, de cada 'subset' reconhecido:
→")
print(corrs)
```

Índices em 'df\_raw' de localização inicial dos 'subsets' reconhecidos:

[365, 2153, 2154, 2155, 9737, 15067, 15200, 15201, 15202, 15203, 15204, 15205, 15206, 19551, 19552, 19553, 20034, 20035, 20036, 20037, 22018, 22019, 22020, 27399, 28115, 28116, 28117, 28118, 28119, 28120]

-----

Coeficientes de correlação, na mesma ordem, de cada 'subset' reconhecido:

[0.9530217193164404, 0.9525798667728674, 0.9574178211979937, 0.9583498564537077, 0.9543742135751833, 0.9526560886216351, 0.9530390541096287, 0.9689221839351155, 0.985407153699156, 1.0, 0.9851166455505797, 0.9679102687040013, 0.951059030176661, 0.9537548707084549, 0.9588287028528437, 0.9560382996174929, 0.9614444666919035, 0.9623493751564967, 0.9626586331918893, 0.9539083314006024, 0.9520593861503935, 0.9542572959099037, 0.9511307978803133, 0.950084255505282, 0.9514522133755656, 0.9594018929324194, 0.9664292094609749, 0.9732573257375167, 0.9712667857760724, 0.9584364247311837]

[21]: *# Mostrando graficamente os 'subsets' reconhecidos...*

```
fig = plt.figure(figsize=(16,70))
quad_y = 2
quad_x = int(len(subsets) / quad_y)
if len(subsets) % quad_y:
    quad_x += 1
for i in range(len(subsets)):
    # Extrai do 'dataset' de dados brutos o 'subset' reconhecido.
    subset = df_raw.loc[subsets[i]:subsets[i]+subset_default.shape[0]-1]
    ax1 = fig.add_subplot(quad_x, quad_y, i + 1)
    ax1.plot(subset.index.values, subset.CLOSE.values)
    ax1.set_title("Subset reconhecido na posição do índice " + str(subsets[i]) +
                  ". Corr=" + str(round(corrs[i],2)) + ".")
plt.show()
```



Observando os gráficos dos 'subsets' reconhecidos, percebemos com facilidade a semelhança deles com o 'subset' padrão. É exatamente isso que queremos. Porém, temos dois problemas a resolver, nos obrigando a limpar a lista de 'subsets' reconhecidos. O primeiro problema é a presença de 'subsets' muito próximos, sendo vizinhos no 'dataset' de dados brutos, sendo considerados redundâncias. Isso pode ser visto entre os gráficos acima, por exemplo, os 'subsets' em sequência de posições 20034, 20035, 20036 e 20037. O segundo problema é a presença do 'subset' exatamente igual ao padrão, visto que nosso 'subset' padrão é uma cópia de parte desse 'dataset'. Isso significa que nosso 'subset' padrão encontrou ele mesmo no 'dataset' de dados brutos, obviamente com coeficiente de correlação igual a '1.0'. Além disso, esse segundo problema vem acompanhando do primeiro.

Para limpá-los a lista de 'subsets' reconhecidos, vamos executar dois passos. Primeiro, vamos remover os subsets vizinhos, deixando aquele que tiver o maior coeficiente de correlação.

[22]: *# Função de remoção de 'subsets' vizinhos...*

```
def subsets_cleaned(subsets, corrs):
    list_subsets_cleaned = []
    list_corrs_cleaned = []
    for i in range(len(subsets)):
        if i > 0:
            subset_current = subsets[i]
            corr_current = corrs[i]
            if subsets[i-1] == subset_current - 1:
                if corr_current > list_corrs_cleaned[-1]:
                    list_subsets_cleaned[-1] = subset_current
                    list_corrs_cleaned[-1] = corr_current
            else:
                list_subsets_cleaned.append(subset_current)
                list_corrs_cleaned.append(corr_current)
        else:
            list_subsets_cleaned.append(subsets[0])
            list_corrs_cleaned.append(corrs[0])
    return list_subsets_cleaned, list_corrs_cleaned
```

[23]: *# Aplicando a função em nosso resultado...*

```
subsets_cleaned, corrs_cleaned = subsets_cleaned(subsets, corrs)
print("Subsets originais:")
print(subsets)
print("Subsets limpos:")
print(subsets_cleaned)
print("-----")
print("Coeficientes originais:")
print(corrs)
```

```
print("Coeficientes limpos:")
print(corrs_cleaned)
```

```
'Subsets' originais:
[365, 2153, 2154, 2155, 9737, 15067, 15200, 15201, 15202, 15203, 15204, 15205,
15206, 19551, 19552, 19553, 20034, 20035, 20036, 20037, 22018, 22019, 22020,
27399, 28115, 28116, 28117, 28118, 28119, 28120]
'Subsets' limpos:
[365, 2155, 9737, 15067, 15203, 19552, 20036, 22019, 27399, 28118]
```

```
-----
Coeficientes originais:
[0.9530217193164404, 0.9525798667728674, 0.9574178211979937, 0.9583498564537077,
0.9543742135751833, 0.9526560886216351, 0.9530390541096287, 0.9689221839351155,
0.985407153699156, 1.0, 0.9851166455505797, 0.9679102687040013,
0.951059030176661, 0.9537548707084549, 0.9588287028528437, 0.9560382996174929,
0.9614444666919035, 0.9623493751564967, 0.9626586331918893, 0.9539083314006024,
0.9520593861503935, 0.9542572959099037, 0.9511307978803133, 0.950084255505282,
0.9514522133755656, 0.9594018929324194, 0.9664292094609749, 0.9732573257375167,
0.9712667857760724, 0.9584364247311837]
Coeficientes limpos:
[0.9530217193164404, 0.9583498564537077, 0.9543742135751833, 0.9526560886216351,
1.0, 0.9588287028528437, 0.9626586331918893, 0.9542572959099037,
0.950084255505282, 0.9732573257375167]
```

Após a limpeza, ficamos com apenas três 'subsets' reconhecidos, no qual um deles ( de coeficiente '1.0') é o nosso 'subset' padrão dentro do 'dataset' de dados brutos. Podemos simplesmente ignorá-lo, mas numa situação real, esperamos que o 'subset' padrão seja escolhido de outra fonte que não seja o 'dataset' de histórico. Por exemplo, podemos extrair em tempo real no mercado financeiro.

Vamos agora rever os gráficos de linha com os 'subsets' limpos.

```
[24]: # Mostrando graficamente os 'subsets' limpos com gráficos de linha...

fig = plt.figure(figsize=(16,20))
quad_y = 2
quad_x = int(len(subsets_cleaned) / quad_y)
if len(subsets_cleaned) % quad_y:
    quad_x += 1
for i in range(len(subsets_cleaned)):
    # Extraí do 'dataset' de dados brutos o 'subset' reconhecido.
    subset = df_raw.loc[subsets_cleaned[i]:subsets_cleaned[i]+subset_default.
    →shape[0]-1]
    ax1 = fig.add_subplot(quad_x, quad_y, i + 1)
    ax1.plot(subset.index.values, subset.CLOSE.values)
    ax1.set_title("Subset reconhecido na posição do índice " +
    →str(subsets_cleaned[i]) +
    " ". Corr=" + str(round(corrs_cleaned[i],2)) + ".")
plt.show()
```



Ficou agora mais claro e conciso o nosso primeiro reconhecimento de 'subsets'. Podemos ver de forma mais realista por gráficos de 'candlestick', observando a parte reconhecida numa parte maior pertencente ao 'dataset' de dados brutos.

[25]: *# Mostrando graficamente os 'subsets' limpos com gráficos de 'candlestick'...*

```

for i in range(len(subsets_cleaned)):
    # Extraí do 'dataset' de dados brutos uma parte maior que o 'subset'
    →reconhecido.
    subset = df_raw.loc[subsets_cleaned[i]-subset_default.shape[0]:
    →subsets_cleaned[i]+2*subset_default.shape[0]-1]
    trace = go.Candlestick(x=subset.index,
                           open=subset['OPEN'],
                           high=subset['HIGH'],
                           low=subset['LOW'],
                           close=subset['CLOSE'])

    # Destacando a parte reconhecida.
    layout = {'title': "Subset reconhecido na posição do índice " +
    →str(subsets_cleaned[i]) +
        ". Corr=" + str(round(corrs_cleaned[i],2)) + ".",
        'yaxis': {'title': 'Pontos'},
        'xaxis': {'rangeslider': {'visible': False}},
        'shapes': [{'x0': subsets_cleaned[i], 'x1':
    →subsets_cleaned[i]+subset_default.shape[0]-1,
                    'y0': 0, 'y1': 1, 'xref': 'x', 'yref': 'paper',
                    'line': {'color': 'rgb(30,30,30)', 'width': 1}}
        ],
        'annotations': [{'x': subsets_cleaned[i], 'y': 0, 'xref': 'x',
        'yref': 'paper', 'showarrow': False, 'xanchor':
    →'left',
                        'text': "'subset' reconhecido"}
        ]
    }

    data = [trace]

fig = go.Figure(data=data,layout=layout)
py.offline.iplot(fig, filename='simple_candlestick_without_range_slider')

```

Subset reconhecido na posição do índice 365. Corr=0.95.



Subset reconhecido na posição do índice 2155. Corr=0.96.



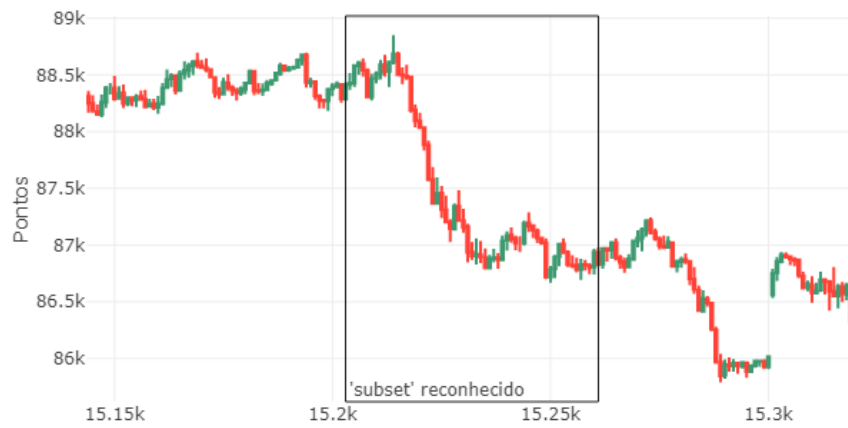
Subset reconhecido na posição do índice 9737. Corr=0.95.



Subset reconhecido na posição do índice 15067. Corr=0.95.



Subset reconhecido na posição do índice 15203. Corr=1.0.



Subset reconhecido na posição do índice 19552. Corr=0.96.



Subset reconhecido na posição do índice 20036. Corr=0.96.



Subset reconhecido na posição do índice 22019. Corr=0.95.



Subset reconhecido na posição do índice 27399. Corr=0.95.



Subset reconhecido na posição do índice 28118. Corr=0.97.

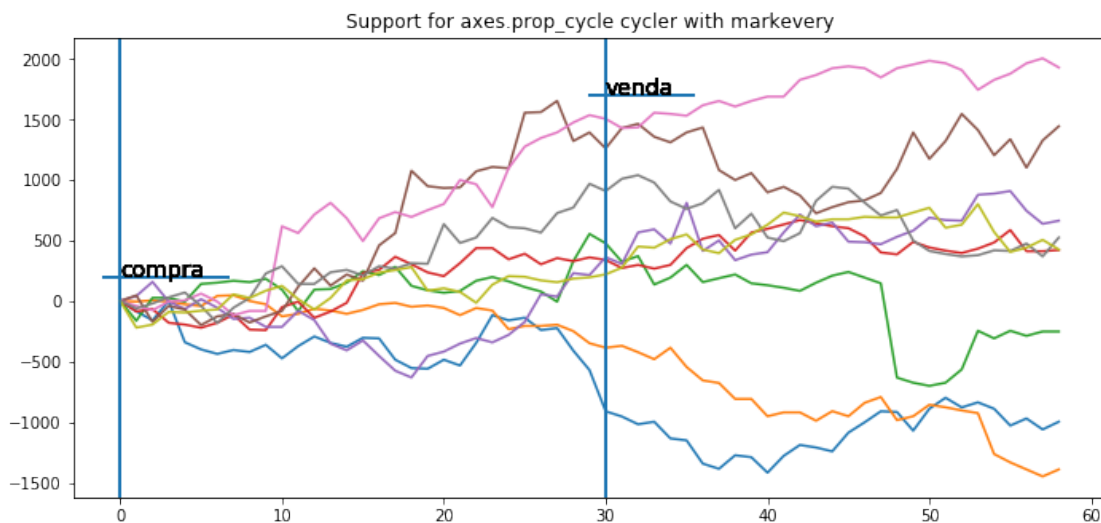


Observando os gráficos de 'candlestick', encontramos uma forma de ajudar o operador de bolsa de valores na sua decisão de comprar lotes desse papel. Supondo que nosso 'subset' padrão foi coletado em tempo real, sendo os últimos 'candles' gerados pela B3 para esse papel, podemos calcular a probabilidade de movimentação da pontuação. Para isso, observamos os nove 'candlesticks' anteriores (desconsiderando o de coeficiente '1.0'), onde em dois deles, após o padrão reconhecido, a pontuação baixou, em um deles se manteve estável e nos seis restante subiu. Se a intenção do operador é fazer uma compra de lotes desse papel, então ele terá neste momento uma probabilidade de aproximadamente 68% de sucesso esperando que a pontuação suba.

A seguir, temos uma visualização de sobreposição dos gráficos de linha após os 'subsets' reconhecidos como forma de visualizar a concentração de movimentação da pontuação.

```
[26]: # Set the plot curve with markers and a title

fig = plt.figure()
ax = fig.add_axes([0.0, 0.0, 1.5, 1.0])
for i in range(len(subsets_cleaned)):
    # Não plota o gráfico do 'subset' padrão.
    if corrs_cleaned[i] != 1.0:
        # Extrai do 'dataset' de dados brutos os 'candles' após os 'subset'
        →reconhecido.
        subset = df_raw.loc[subsets_cleaned[i]+subset_default.shape[0]:
        →subsets_cleaned[i]+2*subset_default.shape[0]-1]
        # Remove a constante incorporada à 'feature' 'CLOSE' subtraindo todos
        →os pontos pelo valor do primeiro ponto.
        subset_CLOSE = [x - subset.CLOSE.values[0] for x in subset.CLOSE]
        # Plota os pontos.
        ax.plot(range(subset_default.shape[0]), subset_CLOSE)
        ax.axvline(x=0, ymin=0, ymax=1)
        ax.axhline(xmin=0.03, y=200, xmax=0.15)
        ax.axvline(x=30, ymin=0, ymax=1)
        ax.axhline(xmin=0.5, y=1700, xmax=0.6)
        ax.text(x=0, y=200, s="compra", fontsize=14)
        ax.text(x=30, y=1700, s="venda", fontsize=14)
plt.title('Support for axes.prop_cycle cyler with markevery')
plt.show()
```



O gráfico acima foi construído com a mesma quantidade de 'candles' dos 'subsets' padrão e reconhecidos, porém, contado a partir do último 'candle' de cada 'subset', mostrando o que viria

depois. Percebemos nesse gráfico que sete dos nove ‘subsets’ reconhecidos teve, em seu histórico, continuidade de movimentação da pontuação de estável para cima até na metade da quantidade de ‘candles’ considerados. A decisão de compra do operador no início faria ele decidir também por vender após metade da quantidade de ‘candles’ do padrão.

Há muitas formas de otimizar a predição da movimentação do preço ou pontuação. Podemos usar, por exemplo, a regressão linear e a própria correlação aplicada à ‘features’ diferentes além da ‘CLOSE’ com ‘CLOSE’ (‘HIGH’ com ‘LOW’, ‘OPEN’ com ‘CLOSE’, ‘FINVOL’ com ‘FINVOL’, ‘TICKVOL’ com ‘TICKVOL’, etc). Por isso, este trabalho significa um ponto de partida para aplicação de novas ferramentas para predição de movimentação de preços ou pontuações de papéis das bolsas e valores.

[ ]: