| NVDA and CAKE Summary Statistics In [import pandas as pd from statsmodels.tsa.stattools import acf, pacf import statsmodels.api as sm import numpy as np import matplotlib.pyplot as plt # trades level schema df = pd.read_csv('xnas-itch-nvidia-cake.csv', parse_dates=['ts_recv', 'ts_event']) # mbo schema: for depth-related questions df_mbo = pd.read_csv('xnas-itch-20240822-mbo.csv', parse_dates=['ts_recv', 'ts_event']) |
|--|
| # split up and prepare dataset # this is returning nothing def prepare_bid_ask_data(df): """ Helper function to create bid_price, ask_price, bid_depth, and ask_depth columns that will be used for later summary statistic calculations.""" # Create bid and ask price columns based on the 'side' column df['bid_price'] = df['price'].where(df['side'] == 'B', None) df['ask_price'] = df['price'].where(df['side'] == 'A', None) # Forward fill the missing bid and ask prices and depths df['bid_price'] = df['bid_price'].ffill() df['ask_price'] = df['ask_price'].ffill() |
| <pre>return df nvda_df = prepare_bid_ask_data(df[df['symbol'] == 'NVDA'].copy()) cake_df = prepare_bid_ask_data(df[df['symbol'] == 'CAKE'].copy()) nvda_df_mbo = df_mbo[df_mbo['symbol'] == 'NVDA'].copy() cake_df_mbo = df_mbo[df_mbo['symbol'] == 'CAKE'].copy() In [nvda_df.columns Out[3]: Index(['ts_recv', 'ts_event', 'rtype', 'publisher_id', 'instrument_id',</pre> |
| Helper functions - contain functionality that calculates results In [import warnings from pandas.errors import SettingWithCopyWarning warnings.simplefilter(action='ignore', category=SettingWithCopyWarning) def calculate_5s_price_impact(df): "" Helper function that contains logic to calculate 5 second price impact."" # Calculate the midpoint price df('midpoint') = (df('bid_price') + df('ask_price')) / 2 # Lag midpoint by 5 seconds df('midpoint_lag') = df('midpoint').shift(5) |
| # Calculate the midpoint return df('midpoint_return') = df('midpoint_lag') # Create a trade sign column based on the 'side' (-1 for 'Ask', +1 for 'Bid') df('trade_sign') = df('side').apply(lambda x: 1 if x == 'B' else -1) # Remove NaN values that result from the shift clean = df.ddropna(subset=['midpoint_return', 'trade_sign']) # Perform the regression of 5-second midpoint return on trade sign price_impact_model = sm.OLS(pd.to_numeric(clean['midpoint_return']), clean['trade_sign']).fit() # Display the summary of the regression model return price_impact_model.summary() |
| <pre>def calculate_depth_at_2x_spread(df_mbo): # Filter for only relevant actions to speed up book calculations df_mbo_filtered = df_mbo[df_mbo['action'].isin(['A', 'C'])].copy() # Group data by second df_mbo_filtered['second'] = df_mbo_filtered['ts_event'].dt.floor('S') # Initialize dictionaries to store order books for bids and asks order_book_bids = {} order_book_asks = {} # Initialize lists to track depth for plotting depth_ask_side = [] depth_bid_side = []</pre> |
| <pre># Initialize variables to track midguote and spread midguotes = [] spreads = [] # Iterate by second print(len(idf_mbo_filtered.groupby('minute'))) i = 0 for second, group in df_mbo_filtered.groupby('minute'): print(i) i+1 # Update the order book based on all actions in the group for _, row in group.iterrows(): if row['action'] == 'A': # Add order if row['side'] == 'B': # Add order order_book_bids[row['order_id']] = {'price': row['price'], 'size': row['size']}</pre> |
| <pre>elif row("side"] == 'A':</pre> |
| <pre>midquotes.append(spread) # Calculate the average spread so far avg_spread = sum(spreads) / len(spreads) if spreads else 0 # Calculate depth at twice the average spread depth_ask = sum[order['size'] for order in order_book_asks.values() if order['price'] <= midquote + 2 * avg_spread]) depth_bid = sum[order['size'] for order in order_book_bids.values() if order['price'] >= midquote - 2 * avg_spread]) depth_ask_side.append(depth_ask) depth_bid_side.append(depth_bid) depth_ask_side_series = pd.Series(depth_ask_side) depth_bid_side_series = pd.Series(depth_bid_side)</pre> |
| <pre>def calculate_bbo_spread_and_depth_per_minute(df_mbo): ''' Calculate the best BBO spread and depth per minute from the provided DataFrame. ''' # Filter the data to include only bid and ask prices bid_data = df_mbo[df_mbo['side'] == 'B'] ask_data = df_mbo[df_mbo['side'] == 'A'] # Group by minute bid_data.loc[:, 'minute'] = bid_data['ts_recv'].dt.floor('T') ask_data.loc[:, 'minute'] = ask_data['ts_recv'].dt.floor('T')</pre> |
| # Get the best bid and ask for each minute best bid = bid_data.groupby('minute').agg({'price': 'max', 'size': 'sum'}).reset_index() best_ask = ask_data.groupby('minute').agg({'price': 'min', 'size': 'sum'}).reset_index() # Rescale prices by dividing by le-9 to convert to dollars best_bid('price') = best_bid('price') / le9 best_ask('price') = best_bsk('price') / le9 # Merge bid and ask data to calculate the spread and depth bbo_data = pd.merge(best_bid, best_ask, on='minute', suffixes=('_bid', '_ask')) # Calculate BBO Spread (Best Ask - Best Bid) bbo_spread = abs(bbo_data['price_ask'] - bbo_data['price_bid']) # Calculate BBO Depth (Sum of Sizes at Best Bid and Best Ask) bbo_depth_bid = bbo_data['size_bid'] bbo_depth_bid = bbo_data['size_bid'] |
| # Set the minute as the index for both series bbo_spread.index = bbo_data['minute'] bbo_depth_bid.index = bbo_data['minute'] bbo_depth_ask.index = bbo_data['minute'] return bbo_spread, bbo_depth_bid, bbo_depth_ask import pandas as pd import numpy as np import statsmodels.api as sm import matplotlib.pyplot as plt def calculate_5s_price_impact_per_minute(df): """ def calculate_5s_price_impact_per_minute(df): """ |
| Helper function to calculate 5-second price impact per minute by regressing the 5-second midpoint quote return on the current trade sign (-1 or +1). ##"" # Calculate the midguote for each row df['midguote'] = (df['price'][df['side'] == 'A']) / 2 # Forward-fill the midguote to ensure we have a midguote available at each second df['midguote'] = df['midguote'].fffill() # Resample the DataFrame to ensure consistent time intervals at the second level df_resampled = df.est_index('ts_event').resample('S').last().ffill() # Remove any duplicate timestamps from the resampled DataFrame df_resampled = df_resampled_ref_resampled_index.duplicated(keep='last')] # Filter the DataFrame to only include trades ('T') |
| <pre>df_trades = df[df['action'] == 'T'].copy() # Assign trade sign: +1 for buy, -1 for sell df_trades['trade_sign'] = np.where(df_trades['side'] == 'B', +1, -1) # Ensure df_resampled has a datetime index for efficient lookups df_resampled = df_resampled.sort_index() # Debugging step: Ensure that there are no duplicates if df_resampled.index.duplicated().any(): raise ValueError("Duplicates detected in the index after processing.") # Get the last available midguote before each trade using asof df_trades['prev_midquote'] = df_resampled('midquote').asof(df_trades['ts_event']) # Efficiently compute the midguote 5 seconds after each trade, using the last available midguote</pre> |
| <pre>df_trades['midquote_5sec'] = df_resampled['midquote'].asof(df_trades['ts_event'] + pd.Timedelta(seconds=5)) # Calculate the 5-second return on the midquote df_trades['midquote_return_5sec'] = (df_trades['midquote_5sec'] - df_trades['prev_midquote']) / df_trades['prev_midquote'] # Drop rows where return or trade sign is NaN df_trades = df_trades.dropna(subset=['midquote_return_5sec', 'trade_sign']) # Check if there's enough data to run the regression if not df_trades.empty:</pre> |
| <pre>plt.plot(df_trades['trade_sign'], model.predict(X), color='red', label='Fitted Line') plt.xlabel('Trade_sign') plt.ylabel('S-Second Midquote Return') plt.title('S-Second Midquote Return vs Trade Sign') plt.legend() plt.show() # Plot the residuals plt.figure(figsize=(10, 6)) plt.scatter(model.predict(X), model.resid, alpha=0.5) plt.satthre(noolel.predict(X), inestyle='') plt.xlabel('Fitted Values') plt.ylabel('Residuals') plt.ylabel('Residuals') plt.title('Residuals') plt.title('Residual Plot') plt.show() else:</pre> |
| <pre>def calculate midguote transaction(df): "'' Helper function to calculate one-second and one-minute midguote and transaction price series.''' # Ensure ts_event is datetime and set it as the index df('ts_event'] = pd.to_datetime(df('ts_event']) df.set_index('ts_event', inplace=True) # Calculate the midpoint price # get the highest bid and lowest ask for each minute # and then calculate the midpoint based on that df('midpoint') = (df('bid_price') + df('ask_price')) / 2</pre> |
| # Resample the data to one-second and one-minute intervals, taking the last value in each interval midquote_lsec = df['midpoint'].resample('18').last() midquote_lmin = df('midpoint'].resample('1T').last() # For transaction prices transaction_lsec = df['price'].resample('18').last() transaction_lmin = df('price'].resample('1T').last() # Drop the first NaN value midquote_lsec.dropna(inplace=True) midquote_lsec.dropna(inplace=True) transaction_lsec.dropna(inplace=True) transaction_lmin.dropna(inplace=True) |
| <pre>return midquote_lsec, midquote_lmin, transaction_lsec, transaction_lmin def calculate_log_returns(midquote_lmin, transaction_lmin): "'' Helper function to calculate log returns.''' midquote_log_return_lmin = np.log(midquote_lmin) - np.log(midquote_lmin.shift(1)) transaction_log_return_lmin = np.log(transaction_lmin) - np.log(transaction_lmin.shift(1)) # Drop NaN values resulting from the shift midquote_log_return_lmin = midquote_log_return_lmin.dropna() transaction_log_return_lmin = transaction_log_return_lmin.dropna() # Display the first few log-returns return midquote_log_return_lmin, transaction_log_return_lmin</pre> |
| <pre>def graph tpm opm(tpm, opm, name): "'' Helper function to break up spike in tpm and opm''' if name == 'NVDA': # Filter for 9 AM to 1 PM opm.between_time('09:00', '13:00') tpm_morning = tpm.between_time('09:00', '13:00') # Filter for 1 PM to 4 PM opm_afternoon = opm.between_time('13:00', '16:00') tpm_afternoon = opm.between_time('13:00', '16:00') # Plot Number of Trades between 9 AM and 1 PM tpm_morning.plot(title=name + ' Number Trades per Minute (9 AM - 1 PM)') plt.ylabel('Time (minute)') plt.ylabel('Number Trades') plt.ylabel('Number Trades') plt.ylabel('Number Trades')</pre> |
| <pre># Plot Number of Trades between 1 PM and 4 PM tpm_afternoon.plot(title=name + ' Number Trades per Minute (1 PM - 4 PM)') plt.xlabel('Time (minute)') plt.show() # Plot Number of Orders between 9 AM and 1 PM opm_morning.plot(title=name + ' Number Orders per Minute (9 AM - 1 PM)') plt.xlabel('Time (minute)') plt.ylabel('Number Orders') plt.show() # Plot Number of Orders between 1 PM and 4 PM opm_afternoon.plot(title=name + ' Number Orders per Minute (1 PM - 4 PM)') plt.xlabel('Time (minute)') plt.xlabel('Time (minute)') plt.xlabel('Time (minute)')</pre> |
| <pre>else: # CAKE - plot as usual # Trades per minute tpm.plot(title = name + ' Number Trades per Minute') plt.xlabel('Time (minute)') plt.ylabel('Number Trades') plt.show() # Orders per minute opm.plot(title=name + ' Number Orders per Minute') plt.xlabel('Time (minute)') plt.ylabel('Number Orders') plt.ylabel('Number Orders') plt.ylabel('Number Orders') plt.show()</pre> def graph_bbo_depth(bbo_depth_bid, bbo_depth_ask, name): ''' Helper function to break up spike in tpm and opm''' |
| <pre>if name == 'NVDA': # Filter for 9 AM to 1 PM bid morning = bbo_depth_bid.between_time('09:00', '13:00') ask_morning = bbo_depth_ask.between_time('09:00', '13:00') # Filter for 1 PM to 4 PM bid_afternoon = bbo_depth_bid.between_time('13:00', '16:00') ask_afternoon = bbo_depth_ask.between_time('13:00', '16:00') # Plot depth between 9 AM and 1 PM bid_morning.plot(title=name + 'BBO_Depth (Bid) (9 AM - 1 PM)') plt.xlabel('Time (minute)') plt.ylabel('Depth') plt.show()</pre> |
| # Plot depth between 1 PM and 4 PM bid afternoon.plot(title=name + ' BBO Depth (Bid) (1 PM - 4 PM)') plt.ylabel('Depth') plt.ylabel('Depth') plt.show() # Plot depth between 9 AM and 1 PM ask_morning.plot(title=name + ' BBO Depth (Ask) (9 AM - 1 PM)') plt.xlabel('Time (minute)') plt.ylabel('Depth') plt.show() # Plot depth between 1 PM and 4 PM ask_afternoon.plot(title=name + ' BBO Depth (Ask) (1 PM - 4 PM)') plt.xlabel('Time (minute)') plt.ylabel('Depth') plt.ylabel('Depth') plt.ylabel('Depth') plt.ylabel('Depth') plt.ylabel('Depth') plt.ylabel('Depth') plt.ylabel('Depth') |
| <pre>else: # CAKE - plot as usual # Trades per minute tpm.plot(title = name + ' Number Trades per Minute') plt.xlabel('Time (minute)') plt.show() # Orders per minute opm.plot(title=name + ' Number Orders per Minute') plt.xlabel('Time (minute)') plt.ylabel('Number Orders') plt.ylabel('Number Orders') </pre> |
| Driver function - calculates all statistics and displays / returns to the user In L from statsmodels.graphics.tsaplots import plot_acf def calc_summary_statistics(df, df_bbo, name): ''' Driver function for (a) - (j)''' # Initialize results table results = {} results_1 = {} results_1 = {} # Get the common minute index minute_index = pd.to_datetime(df['ts_event']).dt.floor('T').unique() # (a) Dollar trading volume per minute df['dollar_vol'] = df('price'] * df('size') |
| <pre>dtw=df.groupby(pd.to_datetime(pd.to_datetime(df['ts_event'])).dt.floor('min'))['dollar_vol'].sum() results['dollar_vol'] = dtv dtv.plot(title=name + ' Dollar Volume Per Minute', logy=True) plt.ylabel('Time (minute)') plt.ylabel('Dollar Trading Volume') plt.show() # (b) Number of trades and number of orders per minute df_mbo('minute') = df_mbo('ts_event').dt.floor('T') # Calculate the number of trades per minute (action 'T') tpm = df_mbo(df_mbo('action') == 'T').groupby('minute').size() # Store the results for trades per minute results['number_trades_per_min'] = tpm</pre> |
| <pre># 'Add' action means act of adding order to order book opm = df_mbo[df_mbo['action'] == '\lambda'].groupby('minute').size() results['number_orders_per_min'] = opm # NVDA had a huge spike at 1pm, we break up so its easier to read - send to helper graph_tpm_opm(tpm, opm, name) # (c) Open, close, high, and low prices stock_open = df.iloc[0]['price'] close = df.iloc[0]['price'] high = df['price'].max() low = df['price'].min() print('\n') print(name + ' Open: ', stock_open)</pre> |
| <pre>print(name + ' Close: ', close) print(name + ' High: ', close, '\n') # (d) VWAP per minute vwap_per_min = df.groupby(pd.to_datetime(df['ts_event']).dt.floor('T')).apply(</pre> |
| bbo_spread, bbo_depth_bid, bbo_depth_ask = calculate_bbo_spread_and_depth_per_minute(df_mbo) results['best_bbo_spread'] = bbo_spread results['best_bbo_depth_bid'] = bbo_depth_bid results['best_bbo_depth_ask'] = bbo_depth_ask graph_bbo_depth(bbo_depth_ask) = bbo_depth_ask graph_bbo_depth(bbo_depth_bid, bbo_depth_ask, name) # (f) Depth at twice that day's average spread # takes a long time to run ask_depth_series, bid_depth_series = calculate_depth_at_2x_spread(df_mbo) results['depth_twice_avg_spread_ask'] = ask_depth_series results['depth_twice_avg_spread_bid'] = bid_depth_series bid_depth_series.plot(title=name + ' Depth at 2x Avg. Spread (Bid)') plt.xlabel('Time_(minutes)') |
| plt.ylabel('Depth') plt.show() ask_depth_series.plot(title=name + 'Depth at 2x Avg. Spread (Ask)') plt.xlabel('Time (minutes)') plt.ylabel('Depth') plt.ylabel('Depth') plt.show() # (g) 5-second price impact per minute - needs to be fixed calculate_5s_price_impact_per_minute(df_mbo.copy()) # (h) Midquote and transaction price series - needs to be fixed midquote_lsec, midquote_lmin, transaction_lsec, transaction_lmin = calculate_midquote_transaction(df.copy()) results_1s['midquote_lsec'] = midquote_lsec results_is['transaction_lsec'] = transaction_lsec |
| results['midquote_lmin'] = midquote_lmin.reindex(minute_index, fill_value=0) results['transaction_lmin'] = transaction_lmin.reindex(minute_index, fill_value=0) plt.plot(midquote_lsec, label='Midquote l-sec') plt.plot(transaction_lsec, label='Transaction l-sec') plt.legend() plt.title(name + ' Midquote and Transaction l-sec') plt.show() plt.plot(midquote_lmin, label='Midquote l-min') plt.plot(transaction_lmin, label='Transaction l-min') plt.legend() plt.title(name + ' Midquote and Transaction l-min') plt.title(name + ' Midquote and Transaction l-min') plt.show() # (i) Log returns - needs to be fixed |
| midquote log_return lmin, transaction_log_return_lmin = calculate_log_returns(midquote_log_return_lmin.reindexx(minute_indexx, fill_value=0) results['transaction_log_return_lmin'] = transaction_log_return_lmin.reindex(minute_index, fill_value=0) plt.plot(results['midquote_log_return_lmin']) plt.title(name + ' Midquote Log Returns l-min') plt.show() plt.plot(results['transaction_log_return_lmin']) plt.title(name + ' Midquote Transaction Returns l-min') plt.title(name + ' Midquote Transaction Returns l-min') plt.show() # (j) Realized variance midquote_realized_variance = (midquote_log_return_lmin ** 2).sum() results['midquote_realized_variance'] = midquote_realized_variance # TODO add transaction realized variance |
| <pre>plt.plot(results['midguote_realized_variance']) plt.title(name + ' Midguote Realized Variance') plt.show() # Convert results to DataFrame min_stats = pd.DataFrame(results) sec_stats = pd.DataFrame(results_Is) # (k) Auto-correlation for midguote and transaction returns midguote_pacf = plot_acf(midguote_log_return_lmin) transaction_pacf = plot_acf(transaction_log_return_lmin)</pre> |
| NVDA Summary Statistics In [calc_summary_statistics(nvda_df, nvda_df_mbo, 'NVDA') NVDA Dollar Volume Per Minute |
| 10 ³ 08 ² 2 ⁰ 08 ² 2 ¹ 08 ² 2 ² 08 ² |
| 200 - 100 - |
| 15000 - 125000 - 15000 - 1500 |
| NVDA Number Orders per Minute (9 AM - 1 PM) 1000 - |
| 09:00 10:00 11:00 12:00 13:00 Time (minute) NVDA Number Orders per Minute (1 PM - 4 PM) 175000 - 125000 - 125000 - 100000 - 100000 - 10000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 10000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 100000 - 100000 - 100000 - 100000 - 100000 - 10000000 - 100000 - 100000 - 100000 - 1000000 - 100000 - 10000000 - 10000000 - 1000000 - 10000000 - 100000000 - 10000000000 |
| 75000 - 25000 - 25000 - 13:00 13:30 14:00 14:30 15:00 15:30 Time (minute) NVDA Open: 128.63 NVDA Close: 125.61 NVDA High: 125.61 NVDA Low: 125.61 |
| NVDA VWAP per Minute 130 129 127 126 125 |
| In [CAKE Summary Statistics In [CAKE Dollar Volume Per Minute CAKE Dollar Volume Per Minute |
| 10 ³ 22 140 ³ |
| Time (minute) CAKE Number Trades per Minute 8000 - 8000 - 2000 - |
| CAKE Number Orders per Minute 175000 - 125000 - 1990 - 100000 |
| 25000 - 25000 - 10:00 11:00 12:00 13:00 14:00 15:00 CAKE Open: 38.76 CAKE Close: 38.7 CAKE High: 38.7 CAKE Liow: 38.7 CAKE Low: 38.7 CAKE Low: 38.7 |
| 39.1 - CAKE VMAP per Minute 39.1 - 38.8 - 38.8 - 38.7 - 22 kg ²⁰ |
| ### #### #### #### #### #### #### #### #### |