```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import t

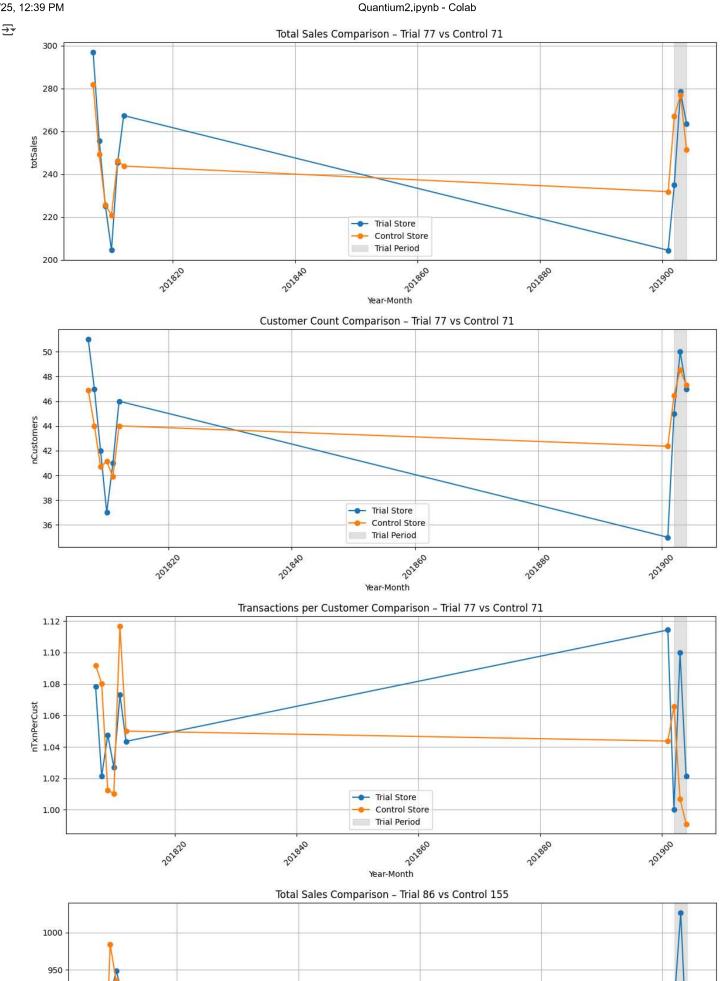
# Step 2: Load and prepare the data
qvi_data = pd.read_csv('/content/QVI_data.csv')
qvi_data['DATE'] = pd.to_datetime(qvi_data['DATE'])
qvi_data['YEARMONTH'] = qvi_data['DATE'].dt.to_period('M').astype(str).str.replace('-', '').astype(int)
print(qvi_data.info())
print(qvi_data.describe())
qvi_data.head()
```

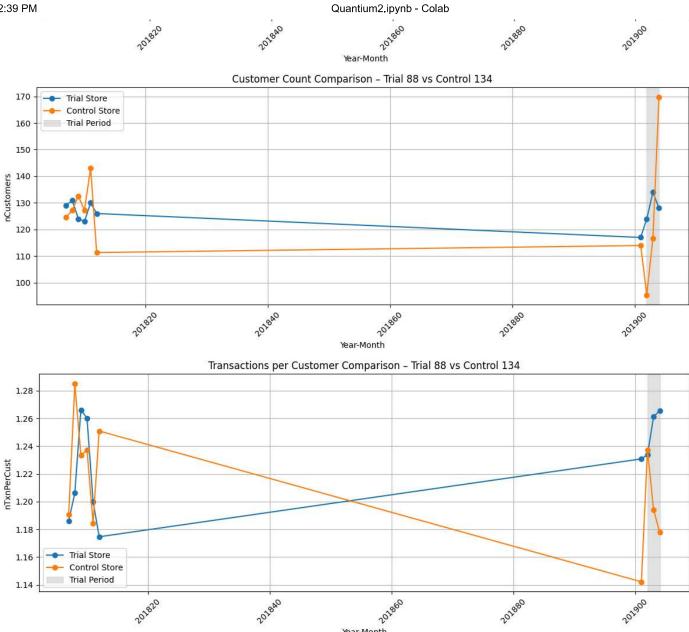
```
→ <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 264834 entries, 0 to 264833
    Data columns (total 13 columns):
         Column
                            Non-Null Count
                                             Dtype
     0
         LYLTY_CARD_NBR
                            264834 non-null
                                              int64
                                             datetime64[ns]
     1
         DATE
                            264834 non-null
     2
         STORE NBR
                            264834 non-null
                                              int64
     3
         TXN_ID
                            264834 non-null
                                              int64
         PROD_NBR
     4
                            264834 non-null
                                              int64
     5
         PROD_NAME
                            264834 non-null
                                              object
         PROD_QTY
                            264834 non-null
         TOT SALES
                            264834 non-null
                                             float64
     8
         PACK_SIZE
                            264834 non-null
                                             int64
     9
         BRAND
                            264834 non-null
                                              object
         LIFESTAGE
                            264834 non-null
                                             object
         PREMTUM CUSTOMER
                            264834 non-null
     11
                                             object
     12
         YEARMONTH
                            264834 non-null
                                             int64
    dtypes: datetime64[ns](1), float64(1), int64(7), object(4)
    memory usage: 26.3+ MB
    None
           LYLTY_CARD_NBR
                                                                STORE_NBR \
                                                      DATE
    count
             2.648340e+05
                                                    264834
                                                            264834.000000
             1.355488e+05
                            2018-12-30 00:52:10.292937984
                                                               135.079423
    mean
    min
             1.000000e+03
                                      2018-07-01 00:00:00
                                                                 1.000000
             7.002100e+04
                                      2018-09-30 00:00:00
                                                                70.000000
    25%
    50%
             1.303570e+05
                                      2018-12-30 00:00:00
                                                               130.000000
                                                               203.000000
    75%
             2.030940e+05
                                      2019-03-31 00:00:00
             2.373711e+06
                                      2019-06-30 00:00:00
                                                               272.000000
    max
             8.057990e+04
                                                                76.784063
    std
                                                       NaN
                  TXN_ID
                               PROD_NBR
                                               PROD_QTY
                                                             TOT_SALES \
           2.648340e+05
                          264834.000000
                                         264834.000000
                                                        264834.000000
    count
           1.351576e+05
                              56.583554
                                               1.905813
                                                              7.299346
    mean
    min
           1.000000e+00
                               1.000000
                                               1.000000
                                                              1.500000
    25%
           6.760050e+04
                              28.000000
                                               2.000000
                                                              5.400000
    50%
           1.351365e+05
                              56.000000
                                               2.000000
                                                              7.400000
    75%
           2.026998e+05
                              85.000000
                                               2.000000
                                                              9.200000
    max
           2.415841e+06
                             114.000000
                                               5.000000
                                                             29.500000
                                               0.343436
    std
           7.813292e+04
                              32.826444
                                                              2.527241
               PACK_SIZE
                               YEARMONTH
           264834.000000
                           264834.000000
    count
              182.425512
                           201856.055163
    mean
    min
               70.000000
                           201807,000000
    25%
              150.000000
                           201809.000000
              170.000000
                           201812.000000
    50%
              175.000000
    75%
                           201903.000000
    max
              380.000000
                           201906.000000
    std
               64.325148
                               47.035278
        LYLTY_CARD_NBR
                         DATE STORE_NBR
                                                                PROD_NAME PROD_QTY TOT_SALES PACK_SIZE
                                                                                                                    BRAND
                                          TXN_ID PROD_NBR
                                                                                                                                    LIFESTAGE
                                                               Natural Chip
                        2018-
                                                                                                                                       YOUNG
     0
                  1000
                                                                  Compny
                                                                                  2
                                                                                            6.0
                                                                                                       175
                                                                                                                 NATURAL
                                                                                                                           SINGLES/COUPLES
                         10-17
                                                               SeaSalt175g
                                                             Red Rock Deli
                        2018-
                                                                                                                                      YOUNG
                  1002
                                               2
                                                              Chikn&Garlic
                                                                                            2.7
                                                         58
                                                                                                       150
                                                                                                                            SINGLES/COUPLES
                        09-16
                                                                 Aioli 150g
                                                               Grain Waves
                        2019-
                                                                     Sour
     2
                  1003
                                               3
                                                                                            3.6
                                                                                                       210
                                                                                                                GRNWVES
                                                                                                                             YOUNG FAMILIES
                                                             Cream&Chives
                        03-07
                                                                     210G
                                                                   Natural
                        2019-
                                                              ChipCo Hony
                  1003
                                                       106
                                                                                                                 NATURAL
                                                                                                                             YOUNG FAMILIES
     3
                                                                                                       175
                                                                                            3.0
                        03-08
                                                                      Sov
                                                                Chckn175g
                                                               WW Original
                        2018-
                                                                                                                                       OLDER
                  1004
                                                             Stacked Chips
                                                                                            1.9
                                                                                                           WOOLWORTHS
                         11-02
                                                                                                                            SINGLES/COUPLES
```

```
# Step 3: Create monthly metrics
monthly_metrics = qvi_data.groupby(['STORE_NBR', 'YEARMONTH']).agg(
   totSales=('TOT_SALES', 'sum'),
   nCustomers=('LYLTY_CARD_NBR', pd.Series.nunique),
   nTransactions=('TXN_ID', pd.Series.nunique),
```

```
nChips=('PROD_QTY', 'sum')
).reset index()
# Step 4: Add derived metrics
monthly\_metrics['nTxnPerCust'] = monthly\_metrics['nTransactions'] \ / \ monthly\_metrics['nCustomers']
monthly_metrics['nChipsPerTxn'] = monthly_metrics['nChips'] / monthly_metrics['nTransactions']
monthly_metrics['avgPricePerUnit'] = monthly_metrics['totSales'] / monthly_metrics['nChips']
# Step 5: Define correlation function
def calculate_correlation(metric_df, store_trial, metric='totSales'):
    trial = metric df[metric df['STORE NBR'] == store trial].sort values('YEARMONTH')
    correlations = []
    for store in metric_df['STORE_NBR'].unique():
        if store == store_trial:
            continue
        control = metric_df[metric_df['STORE_NBR'] == store].sort_values('YEARMONTH')
        merged = pd.merge(trial, control, on='YEARMONTH', suffixes=('_trial', '_ctrl'))
        if merged.empty or len(merged) < 6:</pre>
        corr = merged[f'{metric}_trial'].corr(merged[f'{metric}_ctrl'])
        correlations.append((store trial, store, corr))
    return pd.DataFrame(correlations, columns=['Store1', 'Store2', 'corr measure'])
# Step 6: Define magnitude distance function
def calculate_magnitude(metric_df, store_trial, metric='totSales'):
    trial = metric_df[metric_df['STORE_NBR'] == store_trial]
    distances = []
    for store in metric_df['STORE_NBR'].unique():
        if store == store trial:
            continue
        control = metric_df[metric_df['STORE_NBR'] == store]
        merged = pd.merge(trial, control, on='YEARMONTH', suffixes=('_trial', '_ctrl'))
        if merged.empty or len(merged) < 6:</pre>
            continue
        dist = abs(merged[f'{metric}_trial'] - merged[f'{metric}_ctrl'])
        norm_score = 1 - (dist - dist.min()) / (dist.max() - dist.min() + 1e-9)
        distances.append((store_trial, store, norm_score.mean()))
    return pd.DataFrame(distances, columns=['Store1', 'Store2', 'mag_measure'])
# Step 7: Define control store selector function
def find_best_control_store(trial_store, metric_df, metric='totSales'):
    pretrial_df = metric_df[metric_df['YEARMONTH'] < 201902]</pre>
    corr_df = calculate_correlation(pretrial_df, trial_store, metric=metric)
    mag_df = calculate_magnitude(pretrial_df, trial_store, metric=metric)
    if corr_df.empty or mag_df.empty:
        return None
    control_score = pd.merge(corr_df, mag_df, on=['Store1', 'Store2'])
    control_score['final_score'] = 0.5 * control_score['corr_measure'] + 0.5 * control_score['mag_measure']
    control_score = control_score.sort_values(by='final_score', ascending=False)
    return control score.iloc[0]['Store2'], control score
# Step 8: Apply function to all trial stores
trial_stores = [77, 86, 88]
control_store_results = {}
for trial in trial_stores:
    control_store, full_table = find_best_control_store(trial, monthly_metrics)
    control_store_results[trial] = {
        'control_store': control_store,
        'comparison_table': full_table
    }
for trial, result in control_store_results.items():
    print(f"Trial store {trial} → Best control store: {int(result['control store'])}")
    Trial store 77 → Best control store: 71
     Trial store 86 → Best control store: 155
     Trial store 88 → Best control store: 134
# Function to scale control and plot trial vs control store
def plot_metric_comparison(trial_store, control_store, metric, title=None):
```

```
subset = monthly_metrics[
        (monthly metrics['STORE NBR'].isin([trial store, control store])) &
        (monthly metrics['YEARMONTH'] < 201905)</pre>
    ].copy()
    # Tag store type
    subset['Store_Type'] = subset['STORE_NBR'].apply(lambda x: 'Trial' if x == trial_store else 'Control')
    # Scale control store to match pre-trial levels
    pre_trial = subset['YEARMONTH'] < 201902</pre>
    pre_totals = subset[pre_trial].groupby('Store_Type')[metric].sum()
    scaling_factor = pre_totals['Trial'] / pre_totals['Control']
    subset['scaled_metric'] = subset.apply(
        lambda row: row[metric] * scaling_factor if row['Store_Type'] == 'Control' else row[metric],
    )
    # Plot
    plt.figure(figsize=(12, 5))
    for store_type in ['Trial', 'Control']:
        data = subset[subset['Store_Type'] == store_type]
        plt.plot(data['YEARMONTH'], data['scaled_metric'], marker='o', label=f"{store_type} Store")
    # Highlight trial period
    plt.axvspan(201902, 201904, color='gray', alpha=0.2, label='Trial Period')
    plt.title(title or f"{metric} Comparison - Trial {trial_store} vs Control {control_store}")
    plt.xlabel("Year-Month")
    plt.ylabel(metric)
    plt.xticks(rotation=45)
    plt.grid(True)
    plt.legend()
    plt.tight_layout()
    plt.show()
# Define the metrics you want to visualize
metrics_to_plot = ['totSales', 'nCustomers', 'nTxnPerCust']
metric_titles = {
    'totSales': 'Total Sales Comparison',
    'nCustomers': 'Customer Count Comparison',
    'nTxnPerCust': 'Transactions per Customer Comparison'
}
# Loop through all trial stores and their matched control stores
for trial in trial stores:
    control = int(control_store_results[trial]['control_store'])
    for metric in metrics_to_plot:
        plot_metric_comparison(
            trial_store=trial,
            control store=control,
            metric=metric,
            title=f'{metric_titles[metric]} - Trial {trial} vs Control {control}'
```





Year-Month