

Lab6_Part_1_Time

```
In [3]: import pandas as pd
```

Assignment 1: Date Formats and Date Parts

- First, convert the `date` column to datetime64, by any method.
- Then, create a column representing the time difference between the last date in the data and each date.
- Next, create columns for the date parts year, month, and weekday.
- Finally, format the date to Year-Month-Day (This will be a string/object).

```
In [69]: transactions = pd.read_csv("transactions.csv")
```

```
In [9]: transactions.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
---  --  
 0   date        83488 non-null   object 
 1   store_nbr    83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: int64(2), object(1)
memory usage: 1.9+ MB
```

```
In [13]: # conversion with parse dates in read_csv
```

```
transactions = pd.read_csv("transactions.csv", parse_dates=[ "date"])
transactions
```

Out[13]:

| | date | store_nbr | transactions |
|--------------|------------|-----------|--------------|
| 0 | 2013-01-01 | 25 | 770 |
| 1 | 2013-01-02 | 1 | 2111 |
| 2 | 2013-01-02 | 2 | 2358 |
| 3 | 2013-01-02 | 3 | 3487 |
| 4 | 2013-01-02 | 4 | 1922 |
| ... | ... | ... | ... |
| 83483 | 2017-08-15 | 50 | 2804 |
| 83484 | 2017-08-15 | 51 | 1573 |
| 83485 | 2017-08-15 | 52 | 2255 |
| 83486 | 2017-08-15 | 53 | 932 |
| 83487 | 2017-08-15 | 54 | 802 |

83488 rows × 3 columns

In [15]: `transactions.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   date        83488 non-null   datetime64[ns]
 1   store_nbr   83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: datetime64[ns](1), int64(2)
memory usage: 1.9 MB
```

In [17]: `# conversion with to_datetime
transactions['date'] = pd.to_datetime(transactions['date'])
transactions.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   date        83488 non-null   datetime64[ns]
 1   store_nbr   83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: datetime64[ns](1), int64(2)
memory usage: 1.9 MB
```

In [40]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   date        83488 non-null   datetime64[ns]
 1   store_nbr    83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: datetime64[ns](1), int64(2)
memory usage: 1.9 MB
```

In [19]: *# conversion with astype*

```
transactions['date'] = transactions['date'].astype('datetime64[ns]')
transactions.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   date        83488 non-null   datetime64[ns]
 1   store_nbr    83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: datetime64[ns](1), int64(2)
memory usage: 1.9 MB
```

In [42]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 83488 entries, 0 to 83487
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   date        83488 non-null   datetime64[ns]
 1   store_nbr    83488 non-null   int64  
 2   transactions 83488 non-null   int64  
dtypes: datetime64[ns](1), int64(2)
memory usage: 1.9 MB
```

In [27]: *transactions['date'] = pd.to_datetime(transactions['date'])*
transactions.head()

Out[27]:

| | date | store_nbr | transactions |
|---|------------|-----------|--------------|
| 0 | 2013-01-01 | 25 | 770 |
| 1 | 2013-01-02 | 1 | 2111 |
| 2 | 2013-01-02 | 2 | 2358 |
| 3 | 2013-01-02 | 3 | 3487 |
| 4 | 2013-01-02 | 4 | 1922 |

In [43]:

Out[43]:

| | date | store_nbr | transactions |
|---|------------|-----------|--------------|
| 0 | 2013-01-01 | 25 | 770 |
| 1 | 2013-01-02 | 1 | 2111 |
| 2 | 2013-01-02 | 2 | 2358 |
| 3 | 2013-01-02 | 3 | 3487 |
| 4 | 2013-01-02 | 4 | 1922 |

In [29]:

```
# Calculate the maximum datetime
max_date = transactions['date'].max()
print(max_date)
```

2017-08-15 00:00:00

In []:

In [49]:

```
# Difference between date and max date
transactions['date'] = pd.to_datetime(transactions['date'])
max_date = transactions['date'].max()
transactions['time_to_last_date'] = max_date - transactions['date']

# Dateparts
transactions['year'] = transactions['date'].dt.year
transactions['month'] = transactions['date'].dt.month
transactions['day_of_week'] = transactions['date'].dt.day

# Format Date
transactions['date'] = transactions['date'].dt.strftime('%Y-%B-%d')
transactions[['date', 'store_nbr', 'transactions', 'time_to_last_date', 'year', 'mo
```

Out[49]:

| | date | store_nbr | transactions | time_to_last_date | year | month | day_of_week |
|---|-----------------|-----------|--------------|-------------------|------|-------|-------------|
| 0 | 2013-January-01 | 25 | 770 | 1687 days | 2013 | 1 | 1 |
| 1 | 2013-January-02 | 1 | 2111 | 1686 days | 2013 | 1 | 2 |
| 2 | 2013-January-02 | 2 | 2358 | 1686 days | 2013 | 1 | 2 |
| 3 | 2013-January-02 | 3 | 3487 | 1686 days | 2013 | 1 | 2 |
| 4 | 2013-January-02 | 4 | 1922 | 1686 days | 2013 | 1 | 2 |

Assignment 2: Time Arithmetic

max date in our data was three weeks after 2017-08-15.

- Can you add three weeks to the 'time_to_last_date' column?
- Then, calculate 'weeks_to_last_date' by dividing the number of days in 'time_to_last_date' by 7.

```
In [5]: # overwrite previous transactions df
transactions = pd.read_csv("transactions.csv", parse_dates=["date"])
```

```
In [7]: transactions.tail()
```

```
Out[7]:
```

| | date | store_nbr | transactions |
|--------------|------------|-----------|--------------|
| 83483 | 2017-08-15 | 50 | 2804 |
| 83484 | 2017-08-15 | 51 | 1573 |
| 83485 | 2017-08-15 | 52 | 2255 |
| 83486 | 2017-08-15 | 53 | 932 |
| 83487 | 2017-08-15 | 54 | 802 |

```
In [9]: # recreate columns from assignment 1 using assign
transactions = transactions.assign(
    year=transactions["date"].dt.year,
    month=transactions["date"].dt.month,
    day_of_week=transactions["date"].dt.dayofweek,
    time_to_last_date=transactions["date"].max() - transactions["date"],
)
transactions.head()
```

```
Out[9]:
```

| | date | store_nbr | transactions | year | month | day_of_week | time_to_last_date |
|----------|------------|-----------|--------------|------|-------|-------------|-------------------|
| 0 | 2013-01-01 | 25 | 770 | 2013 | 1 | 1 | 1687 days |
| 1 | 2013-01-02 | 1 | 2111 | 2013 | 1 | 2 | 1686 days |
| 2 | 2013-01-02 | 2 | 2358 | 2013 | 1 | 2 | 1686 days |
| 3 | 2013-01-02 | 3 | 3487 | 2013 | 1 | 2 | 1686 days |
| 4 | 2013-01-02 | 4 | 1922 | 2013 | 1 | 2 | 1686 days |

```
In [11]: # Add three weeks to time to last date column
# Then divide the timedelta (converted to integer) into integer weeks
transactions['time_to_last_date'] = transactions['time_to_last_date'] + pd.Timedelta('3W')
transactions['weeks_to_last_date'] = transactions['time_to_last_date'].dt.days / 7
transactions[['date', 'store_nbr', 'transactions', 'year', 'month', 'day_of_week',
```

Out[11]:

| | date | store_nbr | transactions | year | month | day_of_week | time_to_last_date | weeks_to_l |
|---|------------|-----------|--------------|------|-------|-------------|-------------------|------------|
| 0 | 2013-01-01 | 25 | 770 | 2013 | 1 | 1 | 1708 days | 24 |
| 1 | 2013-01-02 | 1 | 2111 | 2013 | 1 | 2 | 1707 days | 24 |
| 2 | 2013-01-02 | 2 | 2358 | 2013 | 1 | 2 | 1707 days | 24 |
| 3 | 2013-01-02 | 3 | 3487 | 2013 | 1 | 2 | 1707 days | 24 |
| 4 | 2013-01-02 | 4 | 1922 | 2013 | 1 | 2 | 1707 days | 24 |

Assignment 3: Missing Time Series Data

Take a look at the mean value for the oil price using forward fill, backfill, and interpolation. Are they very different?

Then, plot the series with forward fill for:

- The year 2014.
- The month of December 2014.
- The days from December 1st to December 15th, 2014.

In [13]:

```
# Read in oil csv with date as index (and converted to datetime64)
oil = pd.read_csv("oil.csv",
                   index_col="date",
                   parse_dates=True)
```

In [15]:

```
# This is a synonym for datetime64

oil.index.dtype
```

Out[15]:

```
dtype('datetime64[ns]')
```

In [17]:

```
# mean of original series

oil.mean()
```

Out[17]:

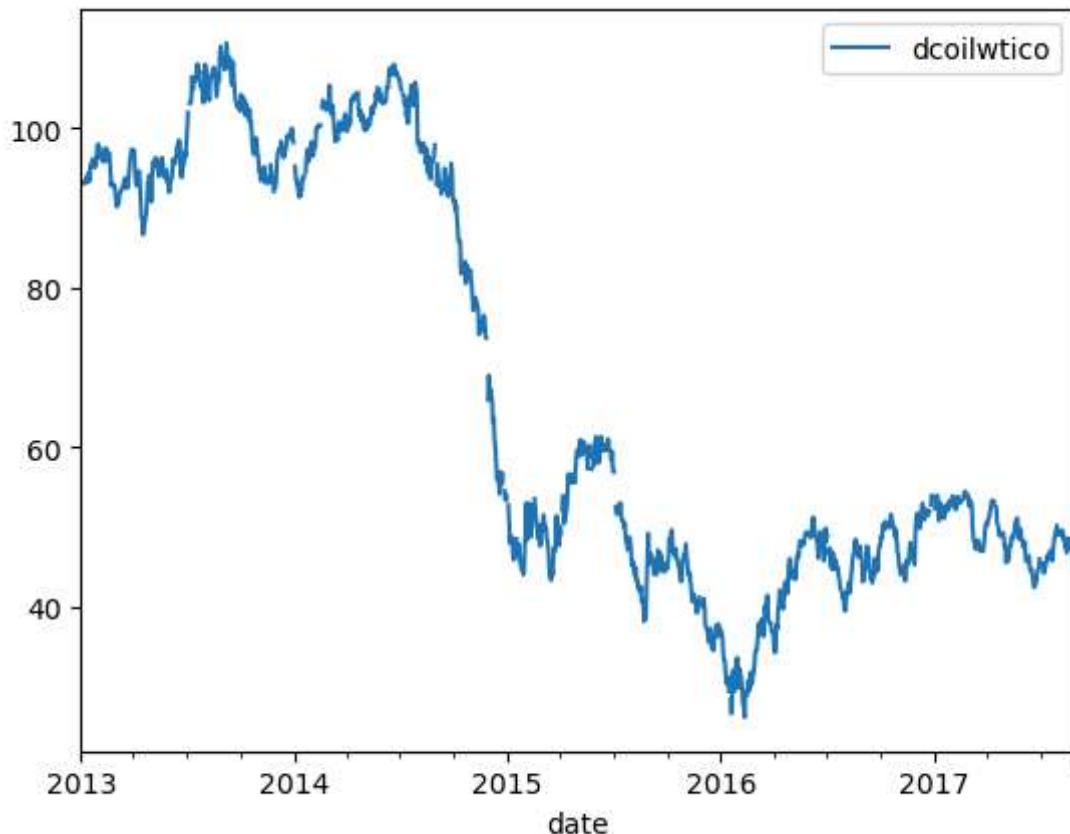
```
dcoilwtico    67.714366
dtype: float64
```

In [19]:

```
# original plot

oil.plot()
```

```
Out[19]: <Axes: xlabel='date'>
```



```
In [21]: # mean of each type of missing value handling for time series
```

```
print(oil.ffill().mean(),
      oil.bfill().mean(),
      oil.interpolate().mean()
     )
```

```
dcoilwtico    67.671249
dtype: float64 dcoilwtico    67.673325
dtype: float64 dcoilwtico    67.661824
dtype: float64
```

```
In [35]: # Filter to 2014 then plot forward filled Series
```

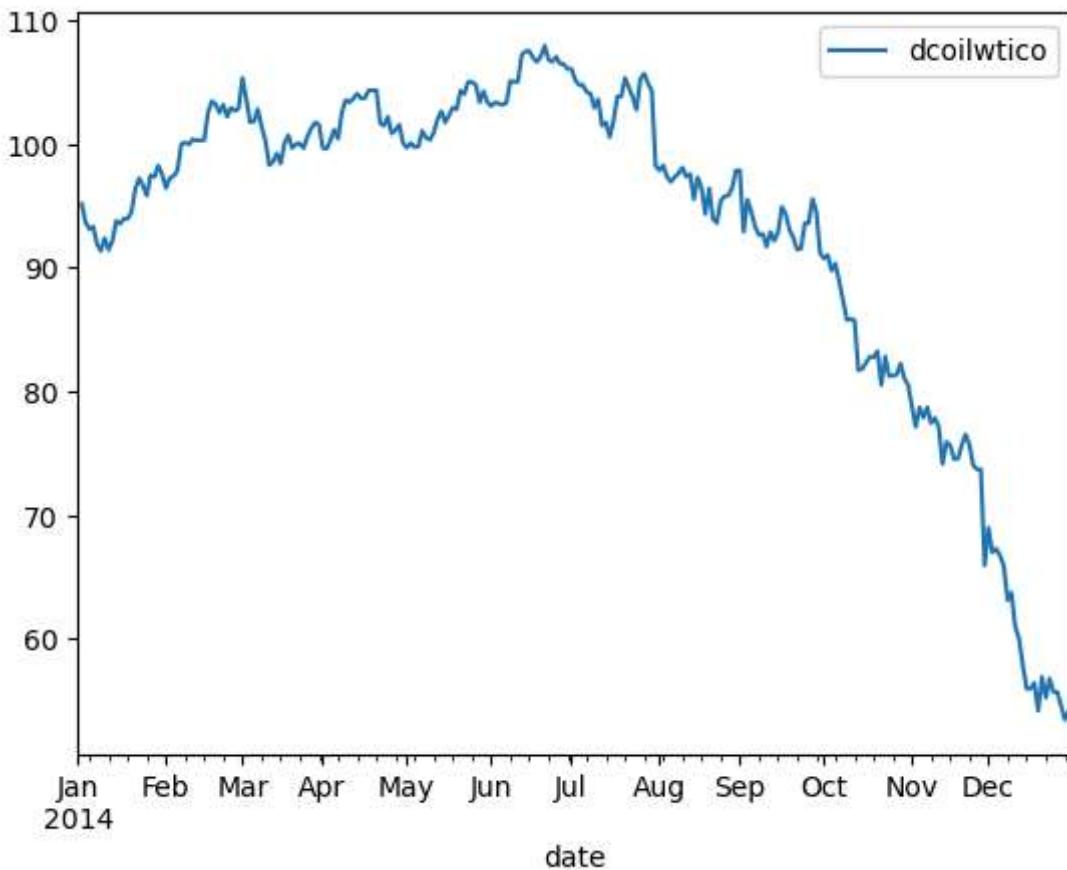
```
oil = pd.read_csv("oil.csv", index_col="date", parse_dates=True)

oil_filtered = oil[oil.index.year == 2014]

oil_filtered = oil_filtered.copy()
oil_filtered['dcoilwtico'] = oil_filtered['dcoilwtico'].ffill()

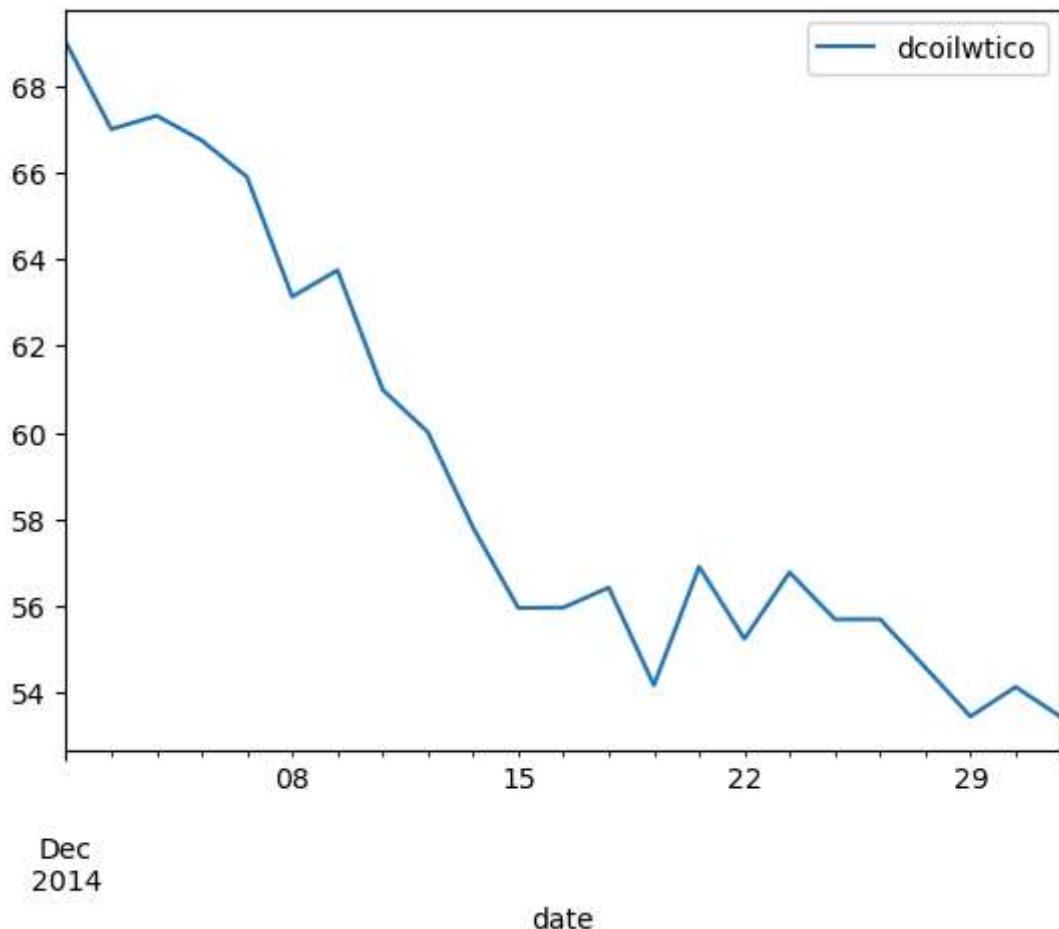
oil_filtered.plot(y='dcoilwtico')
```

```
Out[35]: <Axes: xlabel='date'>
```



```
In [39]: # Filter to December 2014 then plot forward filled Series  
  
oil_filtered = oil[(oil.index.year == 2014) & (oil.index.month == 12)]  
  
oil_filtered = oil_filtered.copy()  
oil_filtered['dcoilwtico'] = oil_filtered['dcoilwtico'].ffill()  
  
oil_filtered.plot(y='dcoilwtico')
```

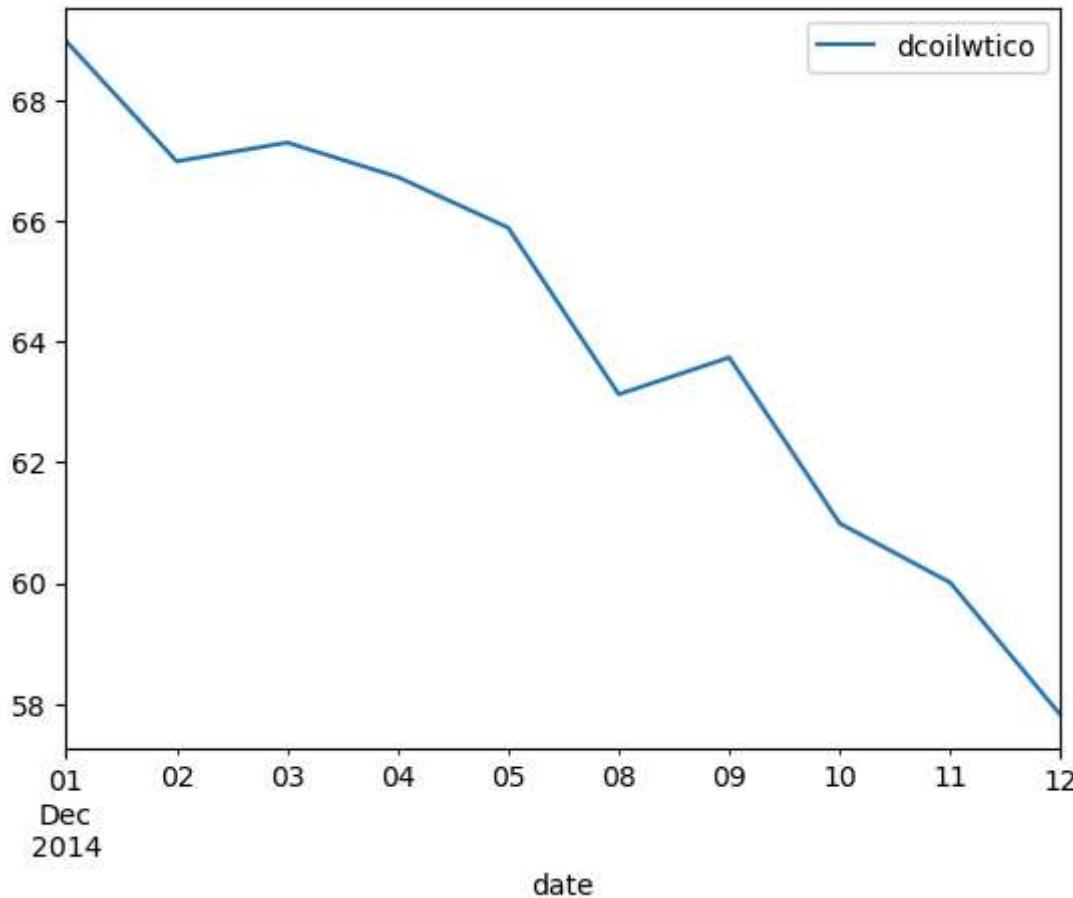
```
Out[39]: <Axes: xlabel='date'>
```



```
In [41]: # Filter to first two weeks of December 2014 then plot forward filled Series
```

```
oil_filtered = oil[(oil.index.year == 2014) & (oil.index.month == 12) & (oil.index.day <= 14)]  
oil_filtered = oil_filtered.copy()  
oil_filtered['dcoilwtico'] = oil_filtered['dcoilwtico'].ffill()  
oil_filtered.plot(y='dcoilwtico')
```

```
Out[41]: <Axes: xlabel='date'>
```



Assignment 4: Shift and Diff

looking into a few different year over year trends related to changes made at store 47.

Can you plot the sum of monthly of transactions in year 2015 vs the sum of monthly transactions in the year prior for store 47?

Make sure to group your DataFrame by year AND month!

Thanks

```
In [45]: # filter df to store 47, 'drop' store_nbr column via Loc
transactions_filtered = transactions[transactions['store_nbr'] == 47].copy()

# Calculate sum of sales by year and month
transactions_filtered['date'] = pd.to_datetime(transactions_filtered['date'])
transactions_filtered.set_index('date', inplace=True)

transactions_filtered['year'] = transactions_filtered.index.year
transactions_filtered['month'] = transactions_filtered.index.month
monthly_sales = transactions_filtered.groupby(['year', 'month'])['transactions'].sum()

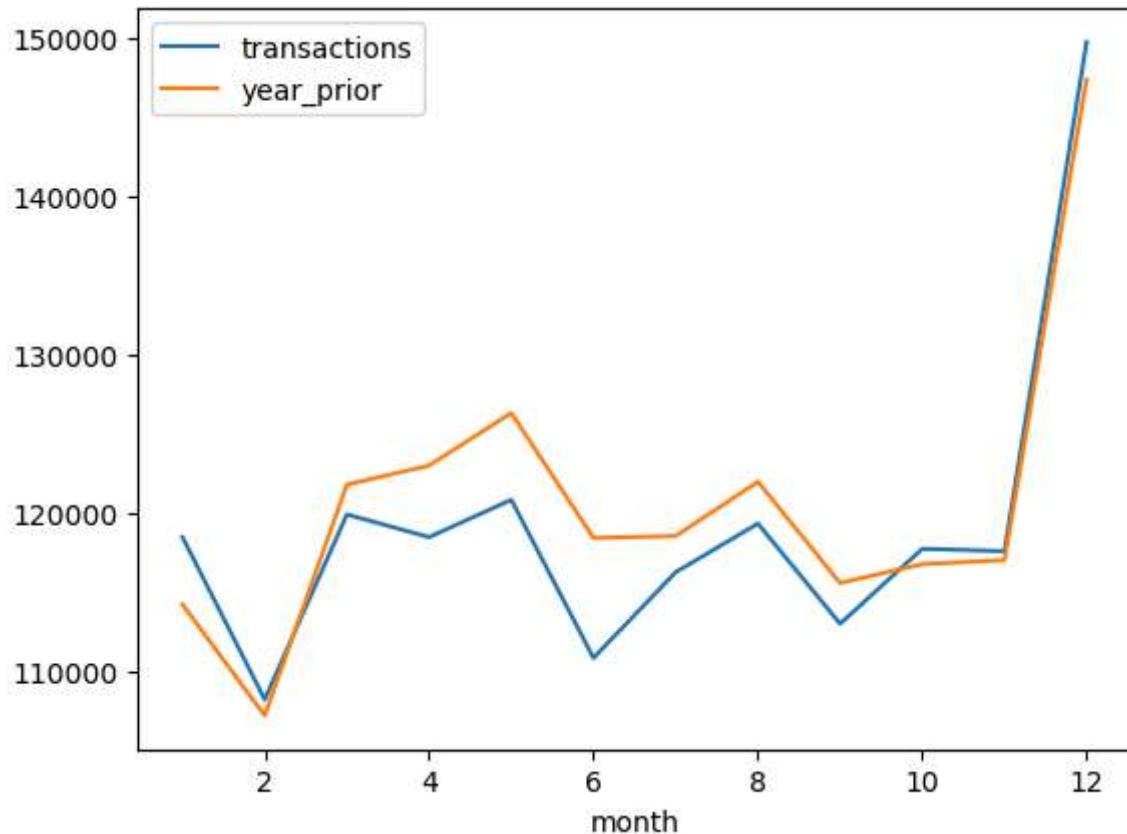
# Calculate a 'year_prior' column by shiftly monthly sales series forward by 12 row
```

```
monthly_sales['year_prior'] = monthly_sales['transactions'].shift(12)

# Filter to 2015 and plot

monthly_sales_2015 = monthly_sales[monthly_sales['year'] == 2015]
monthly_sales_2015.plot(x='month', y=['transactions', 'year_prior'])
```

Out[45]: <Axes: xlabel='month'>



Assignment 5: Resampling Time Series

Plot the monthly and yearly average oil prices.

In [47]: `oil.head()`

Out[47]:

dcoilwtico

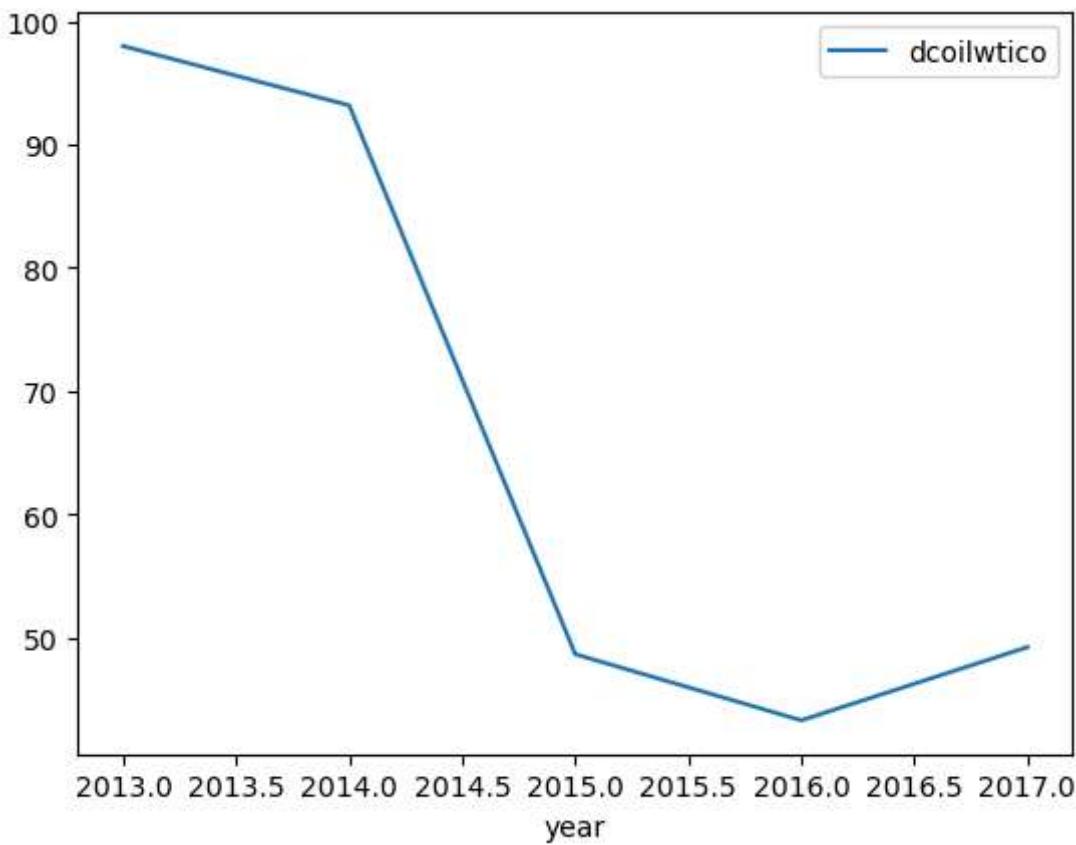
| date | |
|------------|-------|
| 2013-01-01 | NaN |
| 2013-01-02 | 93.14 |
| 2013-01-03 | 92.97 |
| 2013-01-04 | 93.12 |
| 2013-01-07 | 93.20 |

In [57]:

```
# Monthly average oil price
oil['year'] = oil.index.year
yearly_avg_oil = oil.groupby('year')['dcoilwtico'].mean().reset_index()

yearly_avg_oil.plot(x='year', y='dcoilwtico')
```

Out[57]: <Axes: xlabel='year'>



In [67]:

```
# A loop to create various time period averages and plot them
periods = [7, 30, 90]

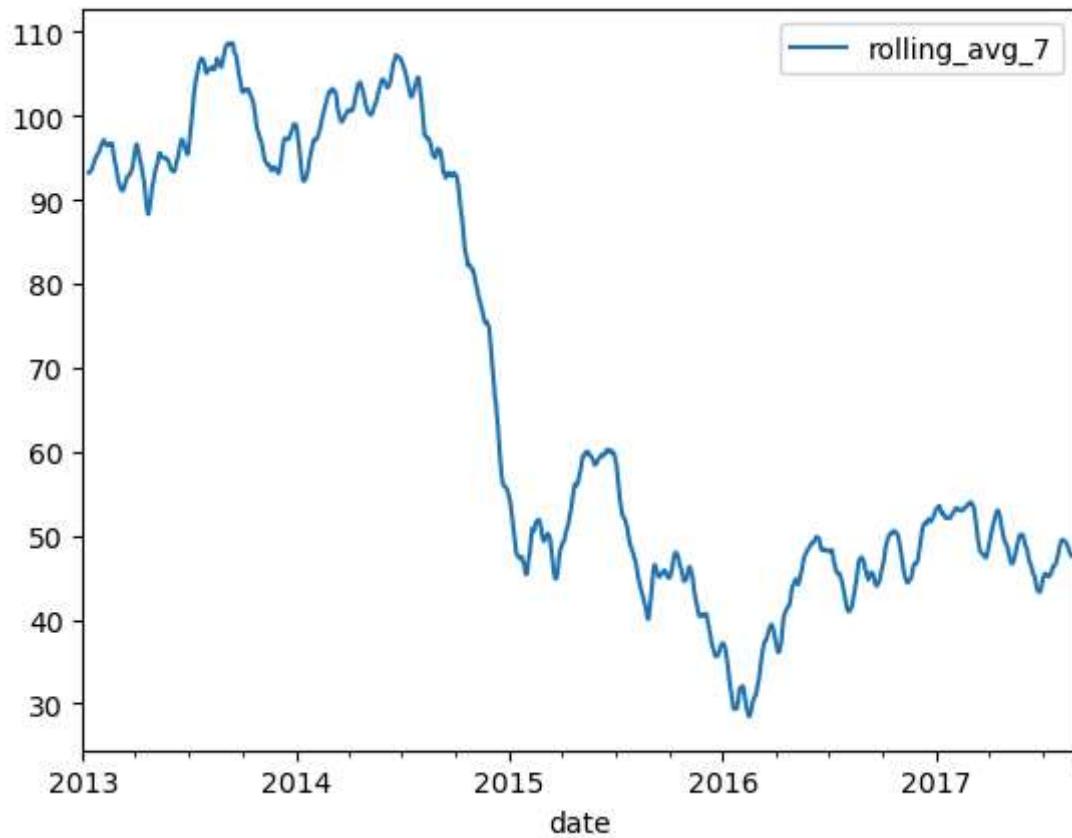
for period in periods:
    oil[f'rolling_avg_{period}'] = oil['dcoilwtico'].rolling(window=period).mean()
    oil[[f'rolling_avg_{period}']].plot()

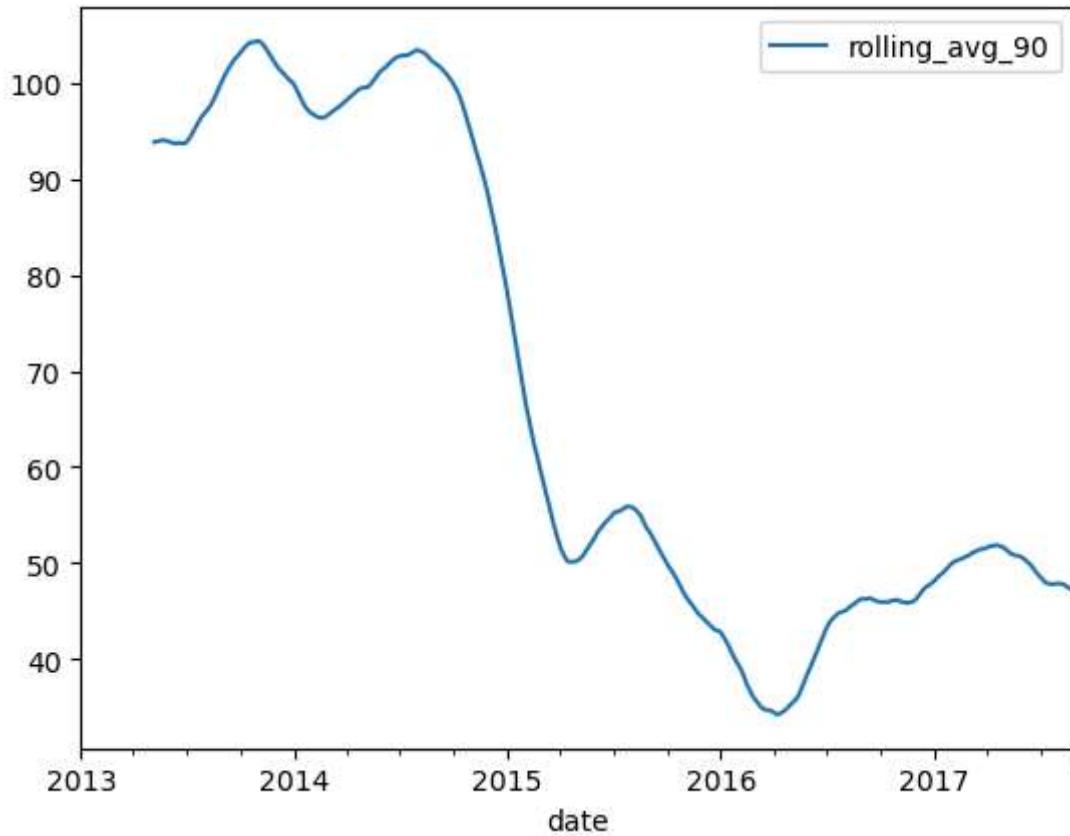
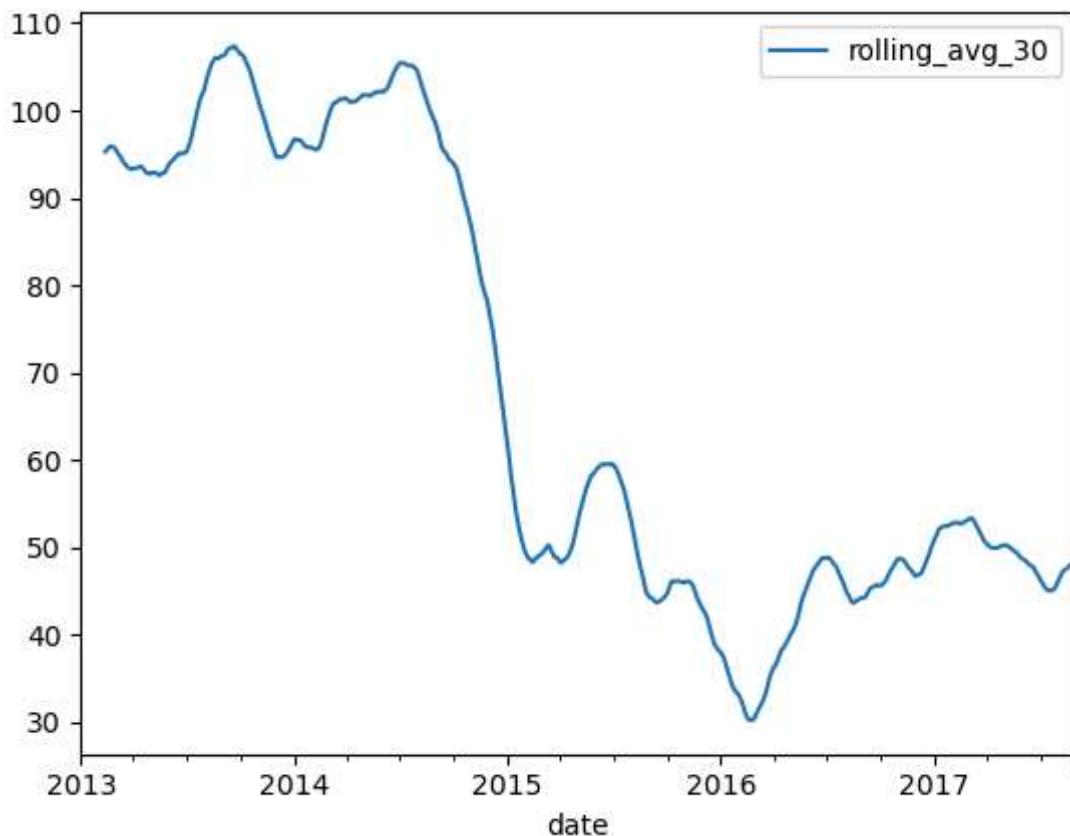
oil['year'] = oil.index.year
```

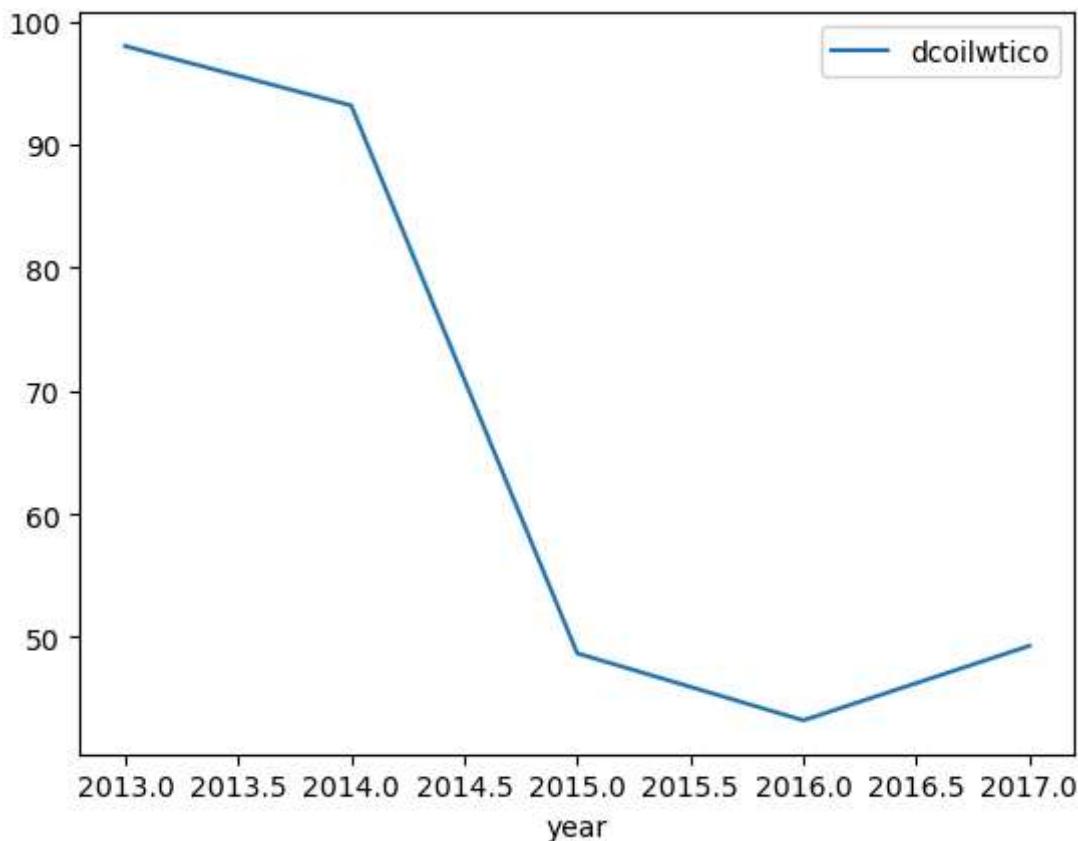
```
yearly_avg_oil = oil.groupby('year')['dcoilwtico'].mean().reset_index()

yearly_avg_oil.plot(x='year', y='dcoilwtico')
```

Out[67]: <Axes: xlabel='year'>







Assignment 6: Rolling Averages

Plot the 90-day moving average for transactions for store 47.

This will help remove some of the noise from our series.

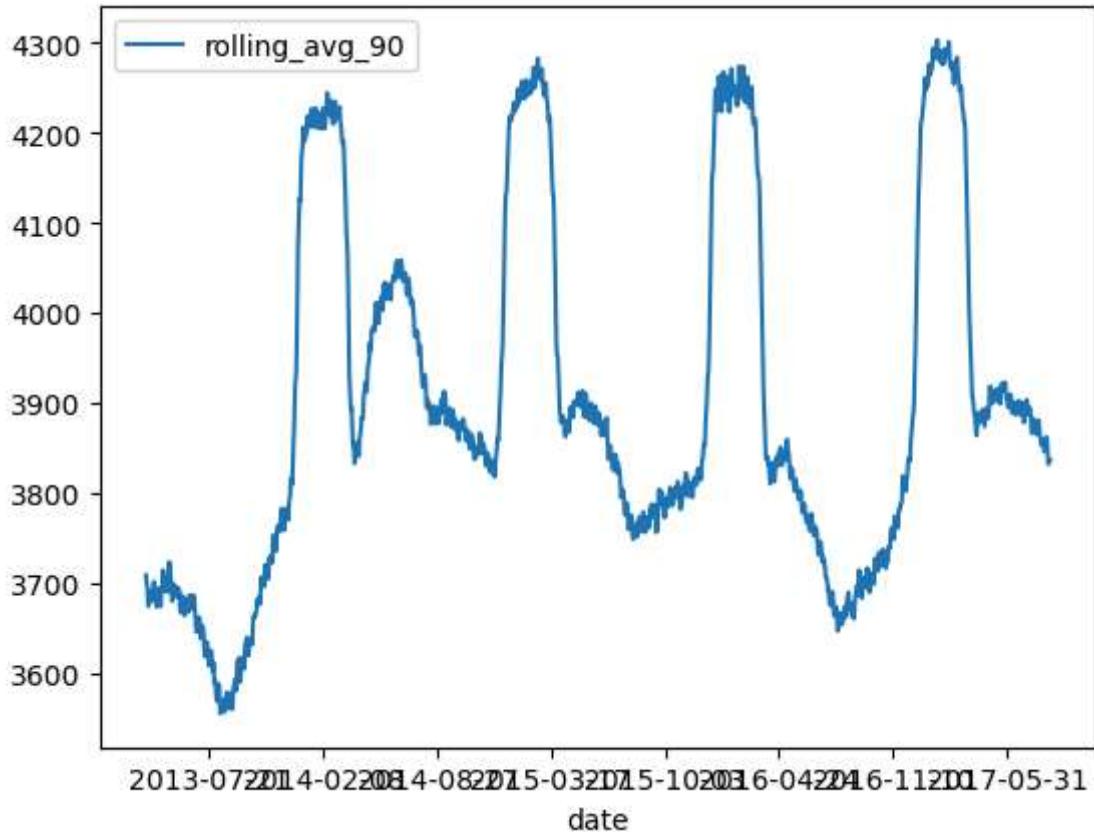
```
In [73]: # recreate transactions47 with date as index
transactions47 = transactions[transactions["store_nbr"]==47].copy()
transactions47.drop(columns=["store_nbr"], inplace=True)
transactions47.set_index("date", inplace=True)
transactions47.head()
```

Out[73]:

| transactions | |
|--------------|------|
| | date |
| 2013-01-02 | 4161 |
| 2013-01-03 | 3660 |
| 2013-01-04 | 3915 |
| 2013-01-05 | 4764 |
| 2013-01-06 | 4935 |

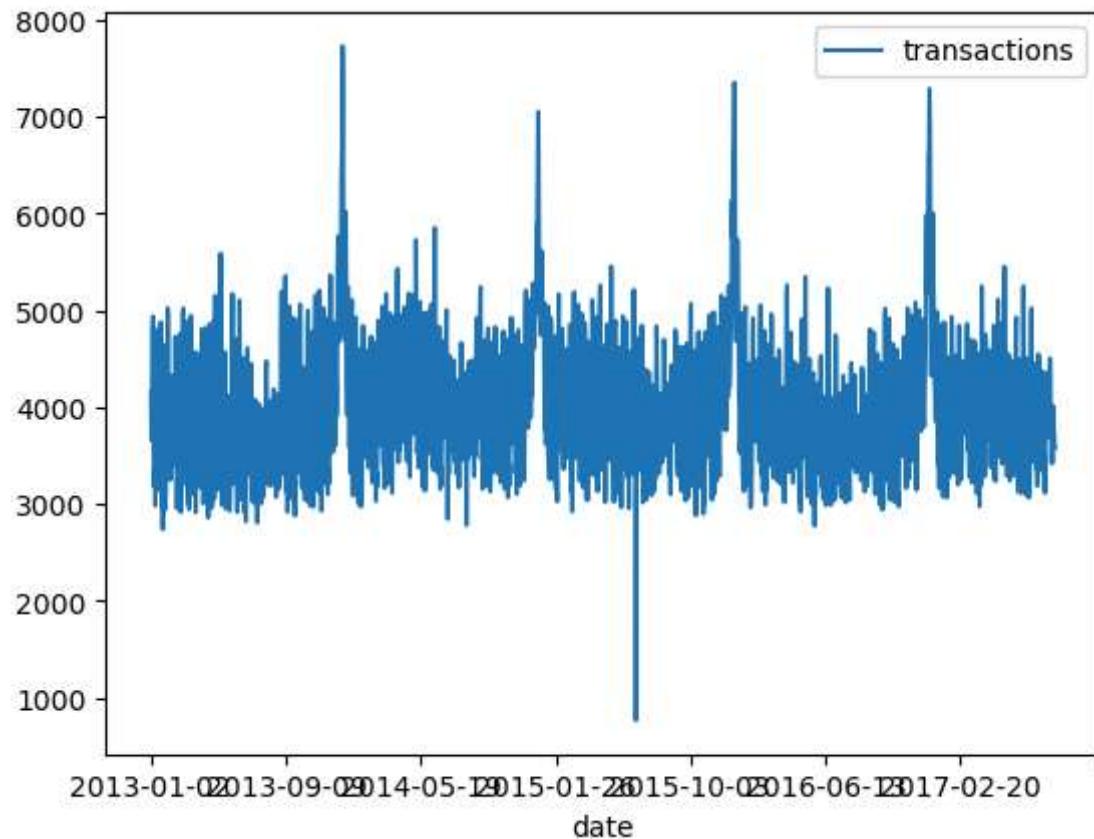
```
In [75]: # Create 90 day rolling average column, drop original transactions column and plot.  
transactions47["rolling_avg_90"] = transactions47["transactions"].rolling(window=90)  
transactions47.drop(columns=["transactions"], inplace=True)  
  
transactions47.plot()
```

Out[75]: <Axes: xlabel='date'>



```
In [83]: # original daily series for comparison  
transactions47 = transactions[transactions["store_nbr"] == 47].set_index("date")[[  
    "transactions"]].plot()
```

Out[83]: <Axes: xlabel='date'>



In []: