# Elements Of Data Science - F2023

# Week 12: Time Series

12/11/2023

# TODOs

- Readings:

  - Recommended: DSFS: Chap 9: Getting Data

  - Recommended: DSFS: Chap 23: Databases and SQL

- HW4, **Due December 15 11:59pm ET**

- Final

  - Review sheet in github repo (soon!)

  - In class, in person.

# Today

- Time Series Transformations

# Questions?

# Environment Setup

# Environment Setup

```python
In [1]:  1  import numpy
         2  import numpy as np
         3  import pandas as pd
         4  import matplotlib.pyplot as plt
         5  import seaborn as sns
         6
         7  sns.set_style('darkgrid')
         8  %matplotlib inline
```

# Time Series

- Data ordered in time

- Applications
  - Financial
  - Economic
  - Scientific
  - etc.

# Time Series Differences

- **Non-i.i.d.** : not independent and identically distributed

- not independent
  - Ex: Stock price

- not-identically distributed
  - Ex: Seasonality

- In other words: Order matters!

# Representing Time in Python

- `datetime` library
- Pandas `Timestamp`

`datetime.date`

# datetime.date

```
In [2]:  1  from datetime import date
         2
         3  friday = date(2022,11,1) # year,month,day
         4  friday
```

Out[2]: datetime.date(2022, 11, 1)

# datetime.date

In [2]:
```python
from datetime import date

friday = date(2022,11,1) # year,month,day
friday
```

Out[2]: datetime.date(2022, 11, 1)

In [3]:
```python
today = date.today()
today
```

Out[3]: datetime.date(2023, 12, 10)

# datetime.date

In [2]:
```python
from datetime import date

friday = date(2022,11,1) # year,month,day
friday
```

Out[2]: datetime.date(2022, 11, 1)

In [3]:
```python
today = date.today()
today
```

Out[3]: datetime.date(2023, 12, 10)

In [4]:
```python
today.year
```

Out[4]: 2023

`datetime.time`

# datetime.time

```
In [5]:  1  from datetime import time
         2
         3  class_start = time(19,10,0) # hour,minute,second,microsecond
         4  class_start

Out[5]: datetime.time(19, 10)
```

# datetime.time

In [5]:
```python
from datetime import time

class_start = time(19,10,0) # hour,minute,second,microsecond
class_start
```

Out[5]: datetime.time(19, 10)

In [6]:
```python
class_start.hour
```

Out[6]: 19

`datetime.datetime`

# datetime.datetime

```
In [7]:  1  from datetime import datetime
         2
         3  # year,month,day,hour,minute,second,microsecond
         4  wednesday_afternoon = datetime(2022,11,30,15)
         5  wednesday_afternoon
```

Out[7]: datetime.datetime(2022, 11, 30, 15, 0)

# datetime.datetime

In [7]:
```python
from datetime import datetime

# year,month,day,hour,minute,second,microsecond
wednesday_afternoon = datetime(2022,11,30,15)
wednesday_afternoon
```

Out[7]: datetime.datetime(2022, 11, 30, 15, 0)

In [8]:
```python
now = datetime.now()
now
```

Out[8]: datetime.datetime(2023, 12, 10, 22, 10, 40, 787053)

`datetime.timedelta`

# datetime.timedelta

```
In [9]: 1 diff = datetime(2022,11,30,1) - datetime(2022,11,29,0)
        2 diff
```

Out[9]: datetime.timedelta(days=1, seconds=3600)

# datetime.timedelta

In [9]:
```
1 diff = datetime(2022,11,30,1) - datetime(2022,11,29,0)
2 diff
```

Out[9]: datetime.timedelta(days=1, seconds=3600)

In [10]:
```
1 diff.total_seconds()
```

Out[10]: 90000.0

# `datetime.timedelta`

In [9]:
```python
1  diff = datetime(2022,11,30,1) - datetime(2022,11,29,0)
2  diff
```

Out[9]: datetime.timedelta(days=1, seconds=3600)

In [10]:
```python
1  diff.total_seconds()
```

Out[10]: 90000.0

In [11]:
```python
1  from datetime import timedelta
2
3  #days,seconds,microseconds,milliseconds,minutes,hours,weeks
4  one_day = timedelta(1)
5
6  date(2022,11,30) + 2*one_day
```

Out[11]: datetime.date(2022, 12, 2)

# Printing Datetimes: `strftime()`

# Printing Datetimes: `strftime()`

```
In [12]:   1  now = datetime.now()
           2  print(now)

           2023-12-10 22:10:40.812133
```

# Printing Datetimes: `strftime()`

```
In [12]:   1  now = datetime.now()
           2  print(now)
```

```
2023-12-10 22:10:40.812133
```

```
In [13]:   1  now.strftime('%a %h %d, %Y %I:%M %p')
```

```
Out[13]:  'Sun Dec 10, 2023 10:10 PM'
```

# Printing Datetimes: `strftime()`

```
In [12]:   1  now = datetime.now()
           2  print(now)
```

2023-12-10 22:10:40.812133

```
In [13]:   1  now.strftime('%a %h %d, %Y %I:%M %p')
```

Out[13]:  'Sun Dec 10, 2023 10:10 PM'

```
1  %Y 4-digit year
2  %y 2-digit year
3  %m 2-digit month
4  %d 2-digit day
5  %H Hour (24-hour)
6  %M 2-digit minute
7  %S 2-digit second
8  ...
```

# Printing Datetimes: `strftime()`

```
In [12]:   1  now = datetime.now()
           2  print(now)
```

```
2023-12-10 22:10:40.812133
```

```
In [13]:   1  now.strftime('%a %h %d, %Y %I:%M %p')
```

```
Out[13]:   'Sun Dec 10, 2023 10:10 PM'
```

```
1 %Y 4-digit year
2 %y 2-digit year
3 %m 2-digit month
4 %d 2-digit day
5 %H Hour (24-hour)
6 %M 2-digit minute
7 %S 2-digit second
8 ...
```

See strftime.org and strfti.me

# Parsing Datetimes: `pandas.to_datetime()`

- `dateutil.parser` available
- pandas has parser built in: `pd.to_datetime()`

# Parsing Datetimes: `pandas.to_datetime()`

- `dateutil.parser` available
- pandas has parser built in: `pd.to_datetime()`

```python
In [14]:   1  pd.to_datetime('11/30/2022 7:36pm')

Out[14]: Timestamp('2022-11-30 19:36:00')
```

# Parsing Datetimes: `pandas.to_datetime()`

- `dateutil.parser` available
- pandas has parser built in: `pd.to_datetime()`

```
In [14]:  1  pd.to_datetime('11/30/2022 7:36pm')

Out[14]:  Timestamp('2022-11-30 19:36:00')
```

```
In [15]:  1  dt_index = pd.to_datetime([datetime(2020, 11, 26),
          2                            '27th of November, 2020',
          3                            '2020-Nov-28',
          4                            '11-29-2030',
          5                            '20201130',
          6                            None
          7                           ])
          8  dt_index

Out[15]:  DatetimeIndex(['2020-11-26', '2020-11-27', '2020-11-28', '2030-11-29',
                         '2020-11-30', 'NaT'],
                        dtype='datetime64[ns]', freq=None)
```

# pandas.Timestamp

- like `datetime.datetime`
- can include **timezone** and **frequency** info
- can handle a missing time: `NaT`
- can be used anywhere `datetime` can be used
- an array of Timestamps can be used as an index

# pandas.Timestamp

- like `datetime.datetime`

- can include **timezone** and **frequency** info

- can handle a missing time: `NaT`

- can be used anywhere `datetime` can be used

- an array of Timestamps can be used as an index

```
In [16]:  1  pd.Timestamp(2022,11,30,19)

Out[16]:  Timestamp('2022-11-30 19:00:00')
```

# `pandas.Timestamp`

- like `datetime.datetime`

- can include **timezone** and **frequency** info

- can handle a missing time: `NaT`

- can be used anywhere `datetime` can be used

- an array of Timestamps can be used as an index

```
In [16]:  1  pd.Timestamp(2022,11,30,19)
```

Out[16]: Timestamp('2022-11-30 19:00:00')

```
In [17]:  1  pd.Timestamp('20221130 7:00pm EST')
```

Out[17]: Timestamp('2022-11-30 19:00:00-0500', tz='tzlocal()')

```
In [18]:  1  pd.Timestamp('20221130 7:00pm',tz='US/Pacific')
```

Out[18]: Timestamp('2022-11-30 19:00:00-0800', tz='US/Pacific')

# **pandas.Timestamp**

- like `datetime.datetime`
- can include **timezone** and **frequency** info
- can handle a missing time: `NaT`
- can be used anywhere `datetime` can be used
- an array of Timestamps can be used as an index

```
In [16]:   1  pd.Timestamp(2022,11,30,19)
```

Out[16]:   Timestamp('2022-11-30 19:00:00')

```
In [17]:   1  pd.Timestamp('20221130 7:00pm EST')
```

Out[17]:   Timestamp('2022-11-30 19:00:00-0500', tz='tzlocal()')

```
In [18]:   1  pd.Timestamp('20221130 7:00pm',tz='US/Pacific')
```

Out[18]:   Timestamp('2022-11-30 19:00:00-0800', tz='US/Pacific')

```
In [19]:   1  dt_index[0]
```

Out[19]:   Timestamp('2020-11-26 00:00:00')

# Accessing Datetime Components with `.dt`

# Accessing Datetime Components with `.dt`

```python
1  df_taxi = pd.read_csv('../data/yellowcab_tripdata_2017-01_subset10000rows.csv',
2                        parse_dates=['tpep_pickup_datetime']).head(3)
3  #df_taxi['tpep_pickup_datetime'] = pd.to_datetime(df_taxi.tpep_pickup_datetime)
4  df_taxi.tpep_pickup_datetime
```

In [20]:

Out[20]:
```
0   2017-01-10 18:37:59
1   2017-01-05 15:14:52
2   2017-01-11 14:47:52
Name: tpep_pickup_datetime, dtype: datetime64[ns]
```

# Accessing Datetime Components with `.dt`

```
In [20]:  1  df_taxi = pd.read_csv('../data/yellowcab_tripdata_2017-01_subset10000rows.csv',
          2                        parse_dates=['tpep_pickup_datetime']).head(3)
          3  #df_taxi['tpep_pickup_datetime'] = pd.to_datetime(df_taxi.tpep_pickup_datetime)
          4  df_taxi.tpep_pickup_datetime
```

```
Out[20]:  0    2017-01-10 18:37:59
          1    2017-01-05 15:14:52
          2    2017-01-11 14:47:52
          Name: tpep_pickup_datetime, dtype: datetime64[ns]
```

```
In [21]:  1  df_taxi.tpep_pickup_datetime.dt.day
```

```
Out[21]:  0    10
          1     5
          2    11
          Name: tpep_pickup_datetime, dtype: int64
```

# Accessing Datetime Components with `.dt`

```python
In [20]:  1  df_taxi = pd.read_csv('../data/yellowcab_tripdata_2017-01_subset10000rows.csv',
          2                        parse_dates=['tpep_pickup_datetime']).head(3)
          3  #df_taxi['tpep_pickup_datetime'] = pd.to_datetime(df_taxi.tpep_pickup_datetime)
          4  df_taxi.tpep_pickup_datetime
```

```
Out[20]:  0   2017-01-10 18:37:59
          1   2017-01-05 15:14:52
          2   2017-01-11 14:47:52
          Name: tpep_pickup_datetime, dtype: datetime64[ns]
```

```python
In [21]:  1  df_taxi.tpep_pickup_datetime.dt.day
```

```
Out[21]:  0    10
          1     5
          2    11
          Name: tpep_pickup_datetime, dtype: int64
```

```python
In [22]:  1  df_taxi.tpep_pickup_datetime.dt.day_of_week  # Monday=0 ... Sunday=6
```

```
Out[22]:  0    1
          1    3
          2    2
          Name: tpep_pickup_datetime, dtype: int64
```

# Accessing Datetime Components with `.dt`

```
In [20]:   1  df_taxi = pd.read_csv('../data/yellowcab_tripdata_2017-01_subset10000rows.csv',
           2                        parse_dates=['tpep_pickup_datetime']).head(3)
           3  #df_taxi['tpep_pickup_datetime'] = pd.to_datetime(df_taxi.tpep_pickup_datetime)
           4  df_taxi.tpep_pickup_datetime
```

```
Out[20]:  0    2017-01-10 18:37:59
          1    2017-01-05 15:14:52
          2    2017-01-11 14:47:52
          Name: tpep_pickup_datetime, dtype: datetime64[ns]
```

```
In [21]:   1  df_taxi.tpep_pickup_datetime.dt.day
```

```
Out[21]:  0    10
          1     5
          2    11
          Name: tpep_pickup_datetime, dtype: int64
```

```
In [22]:   1  df_taxi.tpep_pickup_datetime.dt.day_of_week  # Monday=0 ... Sunday=6
```

```
Out[22]:  0    1
          1    3
          2    2
          Name: tpep_pickup_datetime, dtype: int64
```

```
In [23]:   1  df_taxi.tpep_pickup_datetime.dt.hour
```

```
Out[23]:  0    18
          1    15
          2    14
          Name: tpep_pickup_datetime, dtype: int64
```

# DateIndex Indexing/Selecting/Slicing

# DateIndex Indexing/Selecting/Slicing

```python
In [24]:  1  s = pd.Series(['Dec 1 2021','Jan 2 2022','Feb 3 2022'],
          2              index=pd.to_datetime(['Dec 1 2021','Jan 2 2022','Feb 3 2022']))
          3  s
```

```
Out[24]: 2021-12-01    Dec 1 2021
         2022-01-02    Jan 2 2022
         2022-02-03    Feb 3 2022
         dtype: object
```

# DateIndex Indexing/Selecting/Slicing

```
In [24]:   1  s = pd.Series(['Dec 1 2021','Jan 2 2022','Feb 3 2022'],
           2              index=pd.to_datetime(['Dec 1 2021','Jan 2 2022','Feb 3 2022']))
           3  s
```

```
Out[24]:  2021-12-01    Dec 1 2021
          2022-01-02    Jan 2 2022
          2022-02-03    Feb 3 2022
          dtype: object
```

```
In [25]:   1  # can index normally using iloc
           2  s.iloc[0:2]
```

```
Out[25]:  2021-12-01    Dec 1 2021
          2022-01-02    Jan 2 2022
          dtype: object
```

# DateIndex Indexing/Selecting/Slicing Cont.

# DateIndex Indexing/Selecting/Slicing Cont.

```
In [26]:  1  # only rows from the year 2022
          2  s.loc['2022']
```

```
Out[26]:  2022-01-02     Jan 2 2022
          2022-02-03     Feb 3 2022
          dtype: object
```

# DateIndex Indexing/Selecting/Slicing Cont.

```
In [26]:  1  # only rows from the year 2022
          2  s.loc['2022']
```

```
Out[26]:  2022-01-02    Jan 2 2022
          2022-02-03    Feb 3 2022
          dtype: object
```

```
In [27]:  1  # only rows from January 2022
          2  s.loc['2022-01']
```

```
Out[27]:  2022-01-02    Jan 2 2022
          dtype: object
```

# DateIndex Indexing/Selecting/Slicing Cont.

```
In [26]:    1  # only rows from the year 2022
            2  s.loc['2022']
```

Out[26]:   2022-01-02    Jan 2 2022
           2022-02-03    Feb 3 2022
           dtype: object

```
In [27]:    1  # only rows from January 2022
            2  s.loc['2022-01']
```

Out[27]:   2022-01-02    Jan 2 2022
           dtype: object

```
In [28]:    1  # only rows between Jan 1st 2021 and Jan 2nd 2022, inclusive
            2  s.loc['01/01/2021':'01/02/2022']
```

Out[28]:   2021-12-01    Dec 1 2021
           2022-01-02    Jan 2 2022
           dtype: object

# Datetimes in DataFrames

# Datetimes in DataFrames

```
In [29]:  1  df = pd.DataFrame([['12/1/2021',101,'A'],
          2                     ['1/1/2022',102,'B']],columns=['col1','col2','col3'])
          3  df['col1'] = pd.to_datetime(df.col1)
          4  df.set_index('col1',drop=True,inplace=True)
          5  df
```

Out[29]:

|            | col2 | col3 |
|------------|------|------|
| **col1**   |      |      |
| **2021-12-01** | 101  | A    |
| **2022-01-01** | 102  | B    |

# Datetimes in DataFrames

In [29]:
```python
df = pd.DataFrame([['12/1/2021',101,'A'],
                   ['1/1/2022',102,'B']],columns=['col1','col2','col3'])
df['col1'] = pd.to_datetime(df.col1)
df.set_index('col1',drop=True,inplace=True)
df
```

Out[29]:

|            | col2 | col3 |
|------------|------|------|
| col1       |      |      |
| 2021-12-01 | 101  | A    |
| 2022-01-01 | 102  | B    |

In [30]:
```python
# only return rows from 2022
df.loc['2022']
```

Out[30]:

|            | col2 | col3 |
|------------|------|------|
| col1       |      |      |
| 2022-01-01 | 102  | B    |

# Timestamp Index: Setting Frequency

# Timestamp Index: Setting Frequency

```
In [31]:  1  s = pd.Series(['Nov 1 2022','Nov 3 2022'],index=pd.to_datetime(['Nov 1 2022','Nov 3 2022']))
          2  s
```

```
Out[31]:  2022-11-01    Nov 1 2022
          2022-11-03    Nov 3 2022
          dtype: object
```

# Timestamp Index: Setting Frequency

In [31]:
```python
s = pd.Series(['Nov 1 2022','Nov 3 2022'],index=pd.to_datetime(['Nov 1 2022','Nov 3 2022']))
s
```

Out[31]:
```
2022-11-01     Nov 1 2022
2022-11-03     Nov 3 2022
dtype: object
```

In [32]:
```python
# Use resample() and asfreq() to set frequency
s.resample('D').asfreq()
```

Out[32]:
```
2022-11-01     Nov 1 2022
2022-11-02            NaN
2022-11-03     Nov 3 2022
Freq: D, dtype: object
```

# Timestamp Index: Setting Frequency

```
In [31]: 1 s = pd.Series(['Nov 1 2022','Nov 3 2022'],index=pd.to_datetime(['Nov 1 2022','Nov 3 2022']))
         2 s
```

```
Out[31]: 2022-11-01    Nov 1 2022
         2022-11-03    Nov 3 2022
         dtype: object
```

```
In [32]: 1 # Use resample() and asfreq() to set frequency
         2 s.resample('D').asfreq()
```

```
Out[32]: 2022-11-01    Nov 1 2022
         2022-11-02           NaN
         2022-11-03    Nov 3 2022
         Freq: D, dtype: object
```

```
In [33]: 1 pd.to_datetime(['Nov 1 2022','Nov 3 2022'])
```

```
Out[33]: DatetimeIndex(['2022-11-01', '2022-11-03'], dtype='datetime64[ns]', freq=None)
```

# Timestamp Index: Setting Frequency

In [31]:
```python
s = pd.Series(['Nov 1 2022','Nov 3 2022'],index=pd.to_datetime(['Nov 1 2022','Nov 3 2022']))
s
```

Out[31]:
```
2022-11-01    Nov 1 2022
2022-11-03    Nov 3 2022
dtype: object
```

In [32]:
```python
# Use resample() and asfreq() to set frequency
s.resample('D').asfreq()
```

Out[32]:
```
2022-11-01    Nov 1 2022
2022-11-02           NaN
2022-11-03    Nov 3 2022
Freq: D, dtype: object
```

In [33]:
```python
pd.to_datetime(['Nov 1 2022','Nov 3 2022'])
```

Out[33]:
```
DatetimeIndex(['2022-11-01', '2022-11-03'], dtype='datetime64[ns]', freq=None)
```

In [34]:
```python
# Use date_range with freq to get a range of dates of a certain frequency
pd.date_range(start='Nov 1 2022',end='Nov 3 2022',freq='D')
```

Out[34]:
```
DatetimeIndex(['2022-11-01', '2022-11-02', '2022-11-03'], dtype='datetime64[ns]', freq='D')
```

```
Sample of Available Frequencies
    B    business day frequency
    D    calendar day frequency
    W    weekly frequency
    M    month end frequency
    BM   business month end frequency
    ...
    Q    quarter end frequency
    BQ   business quarter end frequency
    ...
    Y    year end frequency
    BY   business year end frequency
    ...
    BH      business hour frequency
    H       hourly frequency
    T,min   minutely frequency
    S       secondly frequency
    L,ms    milliseconds
    U,us    microseconds
    N       nanoseconds
```

```
 1  Sample of Available Frequencies
 2     B     business day frequency
 3     D     calendar day frequency
 4     W     weekly frequency
 5     M     month end frequency
 6     BM    business month end frequency
 7     ...
 8     Q     quarter end frequency
 9     BQ    business quarter end frequency
10     ...
11     Y     year end frequency
12     BY    business year end frequency
13     ...
14     BH      business hour frequency
15     H       hourly frequency
16     T,min   minutely frequency
17     S       secondly frequency
18     L,ms    milliseconds
19     U,us    microseconds
20     N       nanoseconds
```

https://pandas.pydata.org/pandas-docs/stable/user_guide/timeseries.html#timeseries-offset-aliases

# Timezones

- Handled by `pytz` library

# Timezones

- Handled by `pytz` library

```
In [35]:  1  import pytz
          2
          3  [x for x in pytz.common_timezones if x.startswith('U')]

Out[35]:  ['US/Alaska',
           'US/Arizona',
           'US/Central',
           'US/Eastern',
           'US/Hawaii',
           'US/Mountain',
           'US/Pacific',
           'UTC']
```

# Timezones

- Handled by `pytz` library

```
In [35]:   1  import pytz
           2
           3  [x for x in pytz.common_timezones if x.startswith('U')]

Out[35]:  ['US/Alaska',
           'US/Arizona',
           'US/Central',
           'US/Eastern',
           'US/Hawaii',
           'US/Mountain',
           'US/Pacific',
           'UTC']
```

UTC: coordinated universal time (EST is 5 hours behind, -5:00)

# Timezones Cont.

# Timezones Cont.

```
In [36]:  1  ts = pd.date_range('11/2/2022 9:30am',periods=2,freq='D')
          2  ts
```

Out[36]: DatetimeIndex(['2022-11-02 09:30:00', '2022-11-03 09:30:00'], dtype='datetime64[ns]', freq='D')

# Timezones Cont.

```
In [36]:  1  ts = pd.date_range('11/2/2022 9:30am',periods=2,freq='D')
          2  ts
```

Out[36]: DatetimeIndex(['2022-11-02 09:30:00', '2022-11-03 09:30:00'], dtype='datetime64[ns]', freq='D')

```
In [37]:  1  # Set timezone using .tz_localize()
          2  ts_est = ts.tz_localize('US/Eastern')
          3  ts_est
```

Out[37]: DatetimeIndex(['2022-11-02 09:30:00-04:00', '2022-11-03 09:30:00-04:00'], dtype='datetime64[ns, US/Eastern]', freq=None)

# Timezones Cont.

```
In [36]:  1  ts = pd.date_range('11/2/2022 9:30am',periods=2,freq='D')
          2  ts
```

Out[36]: DatetimeIndex(['2022-11-02 09:30:00', '2022-11-03 09:30:00'], dtype='datetime64[ns]', freq='D')

```
In [37]:  1  # Set timezone using .tz_localize()
          2  ts_est = ts.tz_localize('US/Eastern')
          3  ts_est
```

Out[37]: DatetimeIndex(['2022-11-02 09:30:00-04:00', '2022-11-03 09:30:00-04:00'], dtype='datetime64[ns, US/Eastern]', freq=None)

```
In [38]:  1  # Change timezones using .tz_convert()
          2  ts_est.tz_convert('UTC')
```

Out[38]: DatetimeIndex(['2022-11-02 13:30:00+00:00', '2022-11-03 13:30:00+00:00'], dtype='datetime64[ns, UTC]', freq=None)

# Timezones Cont.

In [36]:
```python
1  ts = pd.date_range('11/2/2022 9:30am',periods=2,freq='D')
2  ts
```

Out[36]: DatetimeIndex(['2022-11-02 09:30:00', '2022-11-03 09:30:00'], dtype='datetime64[ns]', freq='D')

In [37]:
```python
1  # Set timezone using .tz_localize()
2  ts_est = ts.tz_localize('US/Eastern')
3  ts_est
```

Out[37]: DatetimeIndex(['2022-11-02 09:30:00-04:00', '2022-11-03 09:30:00-04:00'], dtype='datetime64[ns, US/Eastern]', freq=None)

In [38]:
```python
1  # Change timezones using .tz_convert()
2  ts_est.tz_convert('UTC')
```

Out[38]: DatetimeIndex(['2022-11-02 13:30:00+00:00', '2022-11-03 13:30:00+00:00'], dtype='datetime64[ns, UTC]', freq=None)

In [39]:
```python
1  # Can also initilize with timezone set
2  ts = pd.date_range('11/2/2022 9:30am',periods=2,freq='D',tz='US/Eastern')
3  ts
```

Out[39]: DatetimeIndex(['2022-11-02 09:30:00-04:00', '2022-11-03 09:30:00-04:00'], dtype='datetime64[ns, US/Eastern]', freq='D')

# Time Series in Python so far:

- `datetime .date .time .datetime .timedelta`
- format with `.strftime()`
- parse time with `pd.to_datetime()`
- `pandas Timestamp Timedelta DatetimeIndex`
- Indexing with `DatetimeIndex`
- Frequencies
- Timezones

Next: Operations on Time Series data

- Shifting
- Resampling
- Moving Windows

# Shifting/Lagging

- Moving data backward or forward in time (lagging/leading)
- Ex: calculate percent change

# Shifting/Lagging

- Moving data backward or forward in time (lagging/leading)

- Ex: calculate percent change

```
In [40]:  1  ts = pd.Series([1,2,8],
          2              index=pd.date_range('1/1/2022',periods=3,freq='M')) # Month End frequency (MS: Month Start)
          3  ts

Out[40]:  2022-01-31    1
          2022-02-28    2
          2022-03-31    8
          Freq: M, dtype: int64
```

# Shifting/Lagging

- Moving data backward or forward in time (lagging/leading)
- Ex: calculate percent change

```
In [40]:  1  ts = pd.Series([1,2,8],
          2              index=pd.date_range('1/1/2022',periods=3,freq='M')) # Month End frequency (MS: Month Start)
          3  ts
```

```
Out[40]: 2022-01-31    1
         2022-02-28    2
         2022-03-31    8
         Freq: M, dtype: int64
```

```
In [41]:  1  ts.shift(1) # last month's value
```

```
Out[41]: 2022-01-31    NaN
         2022-02-28    1.0
         2022-03-31    2.0
         Freq: M, dtype: float64
```

# Shifting

- **percent change**, use one of :
    - (new_value - old_value) / old_value
    - (new_value / old_value) - 1

# Shifting

- **percent change**, use one of :
    - (new_value - old_value) / old_value
    - (new_value / old_value) - 1

```
In [42]:  1  # multiply by 100 to turn into a percent
          2  ((ts / ts.shift(1)) - 1) * 100

Out[42]:  2022-01-31      NaN
          2022-02-28    100.0
          2022-03-31    300.0
          Freq: M, dtype: float64
```

# Example Dataset: Twitter Stock

# Example Dataset: Twitter Stock

```python
In [43]:
1  # from pandas_datareader import data
2  # df_twtr = data.DataReader('TWTR', start='2015', end='11/27/2022', data_source='yahoo')
3  # df_twtr.to_csv('../data/twtr_20150102-20221127.csv')
4  df_twtr = pd.read_csv('../data/twtr_20150102-20221127.csv',parse_dates=['Date'],index_col='Date')
5  df_twtr.head(3).round(2)
```

Out[43]:

|  | High | Low | Open | Close | Volume | Adj Close |
|---|---|---|---|---|---|---|
| **Date** | | | | | | |
| **2015-01-02** | 36.74 | 35.54 | 36.23 | 36.56 | 12062461.0 | 36.56 |
| **2015-01-05** | 37.11 | 35.64 | 36.26 | 36.38 | 15062744.0 | 36.38 |
| **2015-01-06** | 39.45 | 36.04 | 36.27 | 38.76 | 33050812.0 | 38.76 |

# Example Dataset: Twitter Stock

In [43]:
```python
# from pandas_datareader import data
# df_twtr = data.DataReader('TWTR', start='2015', end='11/27/2022', data_source='yahoo')
# df_twtr.to_csv('../data/twtr_20150102-20221127.csv')
df_twtr = pd.read_csv('../data/twtr_20150102-20221127.csv',parse_dates=['Date'],index_col='Date')
df_twtr.head(3).round(2)
```

Out[43]:

| Date | High | Low | Open | Close | Volume | Adj Close |
|------|------|-----|------|-------|--------|-----------|
| 2015-01-02 | 36.74 | 35.54 | 36.23 | 36.56 | 12062461.0 | 36.56 |
| 2015-01-05 | 37.11 | 35.64 | 36.26 | 36.38 | 15062744.0 | 36.38 |
| 2015-01-06 | 39.45 | 36.04 | 36.27 | 38.76 | 33050812.0 | 38.76 |

In [44]:
```python
df_twtr.info() # Adj Close factors in corporate actions, such as stock splits, dividends, and rights offerings
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 1970 entries, 2015-01-02 to 2022-10-27
Data columns (total 6 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   High       1970 non-null   float64
 1   Low        1970 non-null   float64
 2   Open       1970 non-null   float64
 3   Close      1970 non-null   float64
 4   Volume     1970 non-null   float64
 5   Adj Close  1970 non-null   float64
dtypes: float64(6)
memory usage: 107.7 KB
```

# Example Dataset: Twitter Stock

# Example Dataset: Twitter Stock

In [45]:
```python
fig,ax = plt.subplots(1,1,figsize=(24,10))
df_twtr[['Close']].plot(ax=ax,marker='x');
ax.set_ylabel('stock price');
```

# Shifting Example: Percent Change Twitter Close

# Shifting Example: Percent Change Twitter Close

```
In [46]:   1 ((df_twtr.Close / df_twtr.Close.shift(1)) - 1).tail(3).round(3) # # (today / yesterday) - 1
```
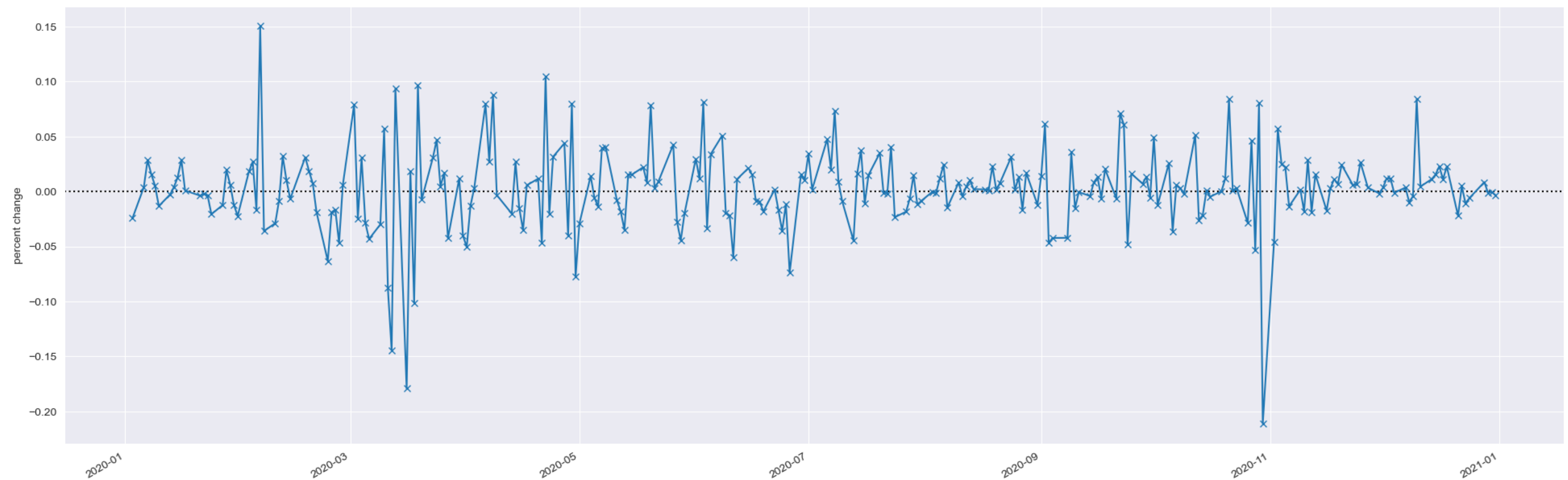
```
Out[46]:   Date
           2022-10-25    0.024
           2022-10-26    0.011
           2022-10-27    0.007
           Name: Close, dtype: float64
```

# Shifting Example: Percent Change Twitter Close

```
In [46]:  1 ((df_twtr.Close / df_twtr.Close.shift(1)) - 1).tail(3).round(3) # # (today / yesterday) - 1
```

```
Out[46]:  Date
          2022-10-25    0.024
          2022-10-26    0.011
          2022-10-27    0.007
          Name: Close, dtype: float64
```

```
In [47]:  1 # plot percent change of close in 2022
          2 fig,ax = plt.subplots(1,1,figsize=(24,8))
          3 close_2020 = df_twtr.loc['2020','Close']
          4 ((close_2020 / close_2020.shift(1)) - 1 ).plot(marker='x',ax=ax,zorder=2);
          5 ax.axhline(ls=':',c='k',zorder=1)
          6 ax.set_ylabel('percent change');
```

# Resampling

- Convert from one frequency to another

- **Downsampling**

  - from higher to lower (day to month)
  - need to aggregate

- **Upsampling**

  - from lower to higher (month to day)
  - need to fill missing

- **Can also be used to set frequency from None**

# Resampling: Initialize Frequency

# Resampling: Initialize Frequency

```
In [48]:    1  df_twtr.index
```

```
Out[48]: DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06', '2015-01-07',
                        '2015-01-08', '2015-01-09', '2015-01-12', '2015-01-13',
                        '2015-01-14', '2015-01-15',
                        ...
                        '2022-10-14', '2022-10-17', '2022-10-18', '2022-10-19',
                        '2022-10-20', '2022-10-21', '2022-10-24', '2022-10-25',
                        '2022-10-26', '2022-10-27'],
                       dtype='datetime64[ns]', name='Date', length=1970, freq=None)
```

# Resampling: Initialize Frequency

```
In [48]:   1  df_twtr.index
```

```
Out[48]:  DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06', '2015-01-07',
                         '2015-01-08', '2015-01-09', '2015-01-12', '2015-01-13',
                         '2015-01-14', '2015-01-15',
                         ...
                         '2022-10-14', '2022-10-17', '2022-10-18', '2022-10-19',
                         '2022-10-20', '2022-10-21', '2022-10-24', '2022-10-25',
                         '2022-10-26', '2022-10-27'],
                        dtype='datetime64[ns]', name='Date', length=1970, freq=None)
```

```
In [49]:   1  df_twtr_B = df_twtr.resample('B').asfreq() # set frequency to business day
           2  df_twtr_B.index
```

```
Out[49]:  DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06', '2015-01-07',
                         '2015-01-08', '2015-01-09', '2015-01-12', '2015-01-13',
                         '2015-01-14', '2015-01-15',
                         ...
                         '2022-10-14', '2022-10-17', '2022-10-18', '2022-10-19',
                         '2022-10-20', '2022-10-21', '2022-10-24', '2022-10-25',
                         '2022-10-26', '2022-10-27'],
                        dtype='datetime64[ns]', name='Date', length=2040, freq='B')
```

# Resampling: Downsampling

- Go from higher/shorter to lower/longer
- Need to aggregate (like groupby)
- Example: Downsampling from business day to business quarter

# Resampling: Downsampling

- Go from higher/shorter to lower/longer

- Need to aggregate (like groupby)

- Example: Downsampling from business day to business quarter

```
In [50]:  1  df_twtr_BQ = df_twtr_B.resample('BQ')
          2  df_twtr_BQ

Out[50]:  <pandas.core.resample.DatetimeIndexResampler object at 0x7fdef51cc250>
```

# Resampling: Downsampling

- Go from higher/shorter to lower/longer

- Need to aggregate (like groupby)

- Example: Downsampling from business day to business quarter

```
In [50]:   1  df_twtr_BQ = df_twtr_B.resample('BQ')
           2  df_twtr_BQ
```

Out[50]: `<pandas.core.resample.DatetimeIndexResampler object at 0x7fdef51cc250>`

```
In [51]:   1  print(df_twtr_BQ)
```

DatetimeIndexResampler [freq=<BusinessQuarterEnd: startingMonth=12>, axis=0, closed=right, label=right, convention=start, origin=start_day]

# Resampling: Downsampling

- Go from higher/shorter to lower/longer

- Need to aggregate (like groupby)

- Example: Downsampling from business day to business quarter

```
In [50]:    1  df_twtr_BQ = df_twtr_B.resample('BQ')
            2  df_twtr_BQ
```

Out[50]: `<pandas.core.resample.DatetimeIndexResampler object at 0x7fdef51cc250>`

```
In [51]:    1  print(df_twtr_BQ)
```

DatetimeIndexResampler [freq=<BusinessQuarterEnd: startingMonth=12>, axis=0, closed=right, label=right, convention=start, origin=start_day]

```
In [52]:    1  df_twtr_BQ.mean().head(3).round(2)
```

Out[52]:

|  | High | Low | Open | Close | Volume | Adj Close |
|---|---|---|---|---|---|---|
| **Date** |  |  |  |  |  |  |
| **2015-03-31** | 45.10 | 43.55 | 44.23 | 44.34 | 20840997.51 | 44.34 |
| **2015-06-30** | 41.63 | 40.38 | 41.17 | 40.87 | 22287099.56 | 40.87 |
| **2015-09-30** | 30.64 | 29.42 | 30.05 | 30.00 | 20065038.11 | 30.00 |

# Resampling: Downsampling

# Resampling: Downsampling

```
In [53]:  1  fig,ax = plt.subplots(1,1,figsize=(24,8))
          2  df_twtr_B.Close.plot(style='-', label='by B',ax=ax)
          3  df_twtr_BQ.Close.mean().plot(style='--',marker='x',label='by BQ',ax=ax)
          4  plt.legend(loc='upper right');
```

# Resampling: Upsampling

- Go from lower/longer to higher/shorter

- Need to decide how to handle missing values

- Example: Upsample from business day to hour

# Resampling: Upsampling

- Go from lower/longer to higher/shorter

- Need to decide how to handle missing values

- Example: Upsample from business day to hour

```
In [54]:   1  df_twtr_B.index[:3]

Out[54]: DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06'], dtype='datetime64[ns]', name='Date', freq='B')
```

# Resampling: Upsampling

- Go from lower/longer to higher/shorter

- Need to decide how to handle missing values

- Example: Upsample from business day to hour

```
In [54]:   1  df_twtr_B.index[:3]
```

```
Out[54]:  DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06'], dtype='datetime64[ns]', name='Date', freq='B')
```

```
In [55]:   1  df_twtr_B.Close.resample('H').asfreq().iloc[0:3].round(2)
```

```
Out[55]:  Date
          2015-01-02 00:00:00     36.56
          2015-01-02 01:00:00       NaN
          2015-01-02 02:00:00       NaN
          Freq: H, Name: Close, dtype: float64
```

# Resampling: Upsampling

- Go from lower/longer to higher/shorter

- Need to decide how to handle missing values

- Example: Upsample from business day to hour

```
In [54]:  1  df_twtr_B.index[:3]
```

```
Out[54]:  DatetimeIndex(['2015-01-02', '2015-01-05', '2015-01-06'], dtype='datetime64[ns]', name='Date', freq='B')
```

```
In [55]:  1  df_twtr_B.Close.resample('H').asfreq().iloc[0:3].round(2)
```

```
Out[55]:  Date
          2015-01-02 00:00:00    36.56
          2015-01-02 01:00:00      NaN
          2015-01-02 02:00:00      NaN
          Freq: H, Name: Close, dtype: float64
```

```
In [56]:  1  df_twtr_B.Close.resample('H').asfreq().iloc[70:73].round(2)
```

```
Out[56]:  Date
          2015-01-04 22:00:00      NaN
          2015-01-04 23:00:00      NaN
          2015-01-05 00:00:00    36.38
          Freq: H, Name: Close, dtype: float64
```

# Resampling: Upsampling

- `ffill()` : Forward Fill

# Resampling: Upsampling

- `ffill()` : Forward Fill

```
In [57]:    1  df_twtr_B.Close.resample('H').ffill().head(3).round(2)

Out[57]:  Date
          2015-01-02 00:00:00     36.56
          2015-01-02 01:00:00     36.56
          2015-01-02 02:00:00     36.56
          Freq: H, Name: Close, dtype: float64
```

# Resampling: Upsampling

- `ffill()` : Forward Fill

```
In [57]:  1 df_twtr_B.Close.resample('H').ffill().head(3).round(2)

Out[57]:  Date
          2015-01-02 00:00:00    36.56
          2015-01-02 01:00:00    36.56
          2015-01-02 02:00:00    36.56
          Freq: H, Name: Close, dtype: float64
```

- `bfill()` : Backward Fill

# Resampling: Upsampling

- `ffill()` : Forward Fill

```
In [57]:  1  df_twtr_B.Close.resample('H').ffill().head(3).round(2)
```

```
Out[57]:  Date
          2015-01-02 00:00:00    36.56
          2015-01-02 01:00:00    36.56
          2015-01-02 02:00:00    36.56
          Freq: H, Name: Close, dtype: float64
```

- `bfill()` : Backward Fill

```
In [58]:  1  df_twtr_B.Close.resample('H').bfill().head(3).round(3)
```

```
Out[58]:  Date
          2015-01-02 00:00:00    36.56
          2015-01-02 01:00:00    36.38
          2015-01-02 02:00:00    36.38
          Freq: H, Name: Close, dtype: float64
```

# Moving/Rolling Windows

- Apply function on a fixed window moving across time
- Method of smoothing out the data
- **center** : place values at center of window

# Moving/Rolling Windows

- Apply function on a fixed window moving across time

- Method of smoothing out the data

- **center** : place values at center of window

```
In [59]:   1 df_twtr_B.Close['2020-11-02':'2020-11-06'].round(2)

Out[59]:   Date
           2020-11-02    39.47
           2020-11-03    41.73
           2020-11-04    42.76
           2020-11-05    43.71
           2020-11-06    43.12
           Freq: B, Name: Close, dtype: float64
```

# Moving/Rolling Windows

- Apply function on a fixed window moving across time

- Method of smoothing out the data

- **center** : place values at center of window

```
In [59]:  1 df_twtr_B.Close['2020-11-02':'2020-11-06'].round(2)
```

```
Out[59]: Date
         2020-11-02    39.47
         2020-11-03    41.73
         2020-11-04    42.76
         2020-11-05    43.71
         2020-11-06    43.12
         Freq: B, Name: Close, dtype: float64
```

```
In [60]:  1 rolling_3 = df_twtr_B.Close['2020-11-02':'2020-11-06'].rolling(3, center=True)
          2 rolling_3
```

```
Out[60]: Rolling [window=3,center=True,axis=0,method=single]
```
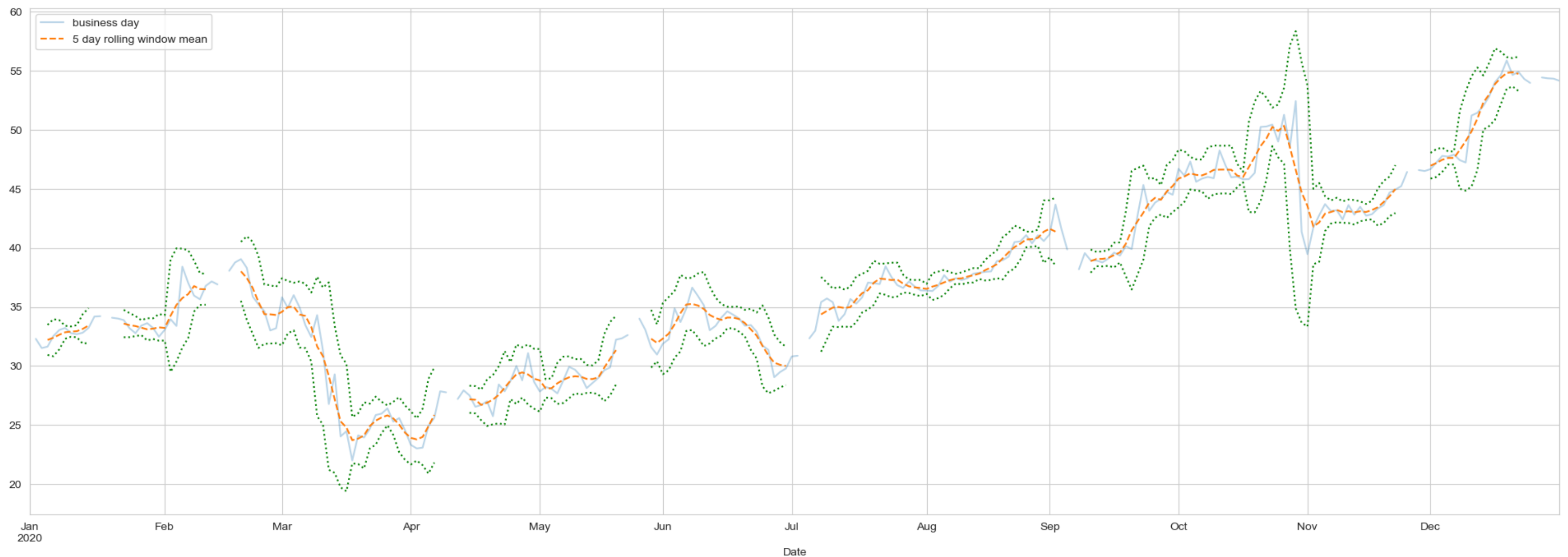
# Moving/Rolling Windows

- Apply function on a fixed window moving across time

- Method of smoothing out the data

- **center** : place values at center of window

```
In [59]:  1  df_twtr_B.Close['2020-11-02':'2020-11-06'].round(2)
```

```
Out[59]:  Date
          2020-11-02    39.47
          2020-11-03    41.73
          2020-11-04    42.76
          2020-11-05    43.71
          2020-11-06    43.12
          Freq: B, Name: Close, dtype: float64
```

```
In [60]:  1  rolling_3 = df_twtr_B.Close['2020-11-02':'2020-11-06'].rolling(3, center=True)
          2  rolling_3
```

```
Out[60]:  Rolling [window=3,center=True,axis=0,method=single]
```

```
In [61]:  1  rolling_3.mean()['2020-11-02':'2020-11-06'].round(2)
```

```
Out[61]:  Date
          2020-11-02     NaN
          2020-11-03    41.32
          2020-11-04    42.73
          2020-11-05    43.20
          2020-11-06     NaN
          Freq: B, Name: Close, dtype: float64
```

# Moving Windows

# Moving Windows

```
1 sns.set_style("whitegrid")
2 rolling = df_twtr_B.Close.rolling(5, center=True)
3
4 fig,ax = plt.subplots(1,1,figsize=(24,8));
5 df_twtr_B.loc['2020'].Close.plot(style='-',alpha=0.3,label='business day');
6 rolling.mean().loc['2020'].plot(style='--',label='5 day rolling window mean');
7 (rolling.mean().loc['2020'] + 2*rolling.std().loc['2020']).plot(style=':',c='g',label='_nolegend_');
8 (rolling.mean().loc['2020'] - 2*rolling.std().loc['2020']).plot(style=':',c='g',label='_nolegend_');
9 ax.legend();
```

# Example: Bike Travel (From PDSH Chapter 3.11)

- Bicycle traffic over Fremont Bridge in Seattle in 2012

- Data gathered using: `!curl -o ../data/FremontBridge.csv https://data.seattle.gov/api/views/65db-xm6k/rows.csv?accessType=DOWNLOAD`

# Example: Bike Travel (From PDSH Chapter 3.11)

- Bicycle traffic over Fremont Bridge in Seattle in 2012

- Data gathered using: `!curl -o ../data/FremontBridge.csv https://data.seattle.gov/api/views/65db-xm6k/rows.csv?accessType=DOWNLOAD`

```
In [63]:   1  df_bike_counts = pd.read_csv('../data/FremontBridge_2012-2015.csv',parse_dates=['Date'],index_col='Date')
           2  df_bike_counts.columns = ['Total','East','West']
           3  df_bike_counts.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 28440 entries, 2012-10-03 00:00:00 to 2015-12-31 23:00:00
Data columns (total 3 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   Total   28433 non-null  float64
 1   East    28433 non-null  float64
 2   West    28433 non-null  float64
dtypes: float64(3)
memory usage: 888.8 KB
```

# Example: Bike Travel (From PDSH Chapter 3.11)

- Bicycle traffic over Fremont Bridge in Seattle in 2012

- Data gathered using: `!curl -o ../data/FremontBridge.csv https://data.seattle.gov/api/views/65db-xm6k/rows.csv?accessType=DOWNLOAD`

```
In [63]:   1 df_bike_counts = pd.read_csv('../data/FremontBridge_2012-2015.csv',parse_dates=['Date'],index_col='Date')
           2 df_bike_counts.columns = ['Total','East','West']
           3 df_bike_counts.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 28440 entries, 2012-10-03 00:00:00 to 2015-12-31 23:00:00
Data columns (total 3 columns):
 #   Column   Non-Null Count   Dtype
---  ------   --------------   -----
 0   Total    28433 non-null   float64
 1   East     28433 non-null   float64
 2   West     28433 non-null   float64
dtypes: float64(3)
memory usage: 888.8 KB
```

```
In [64]:   1 df_bike_counts.head(3)
```

Out[64]:

|                     | Total | East | West |
|---------------------|-------|------|------|
| **Date**            |       |      |      |
| 2012-10-03 00:00:00 | 13.0  | 4.0  | 9.0  |
| 2012-10-03 01:00:00 | 10.0  | 4.0  | 6.0  |
| 2012-10-03 02:00:00 | 2.0   | 1.0  | 1.0  |

# Example: Fill Missing Values

# Example: Fill Missing Values

```
In [65]:  1  f'proportion missing: {sum(df_bike_counts.Total.isna()) / len(df_bike_counts):0.5f}'

Out[65]:  'proportion missing: 0.00025'
```

# Example: Fill Missing Values

```
In [65]:   1   f'proportion missing: {sum(df_bike_counts.Total.isna()) / len(df_bike_counts):0.5f}'
```
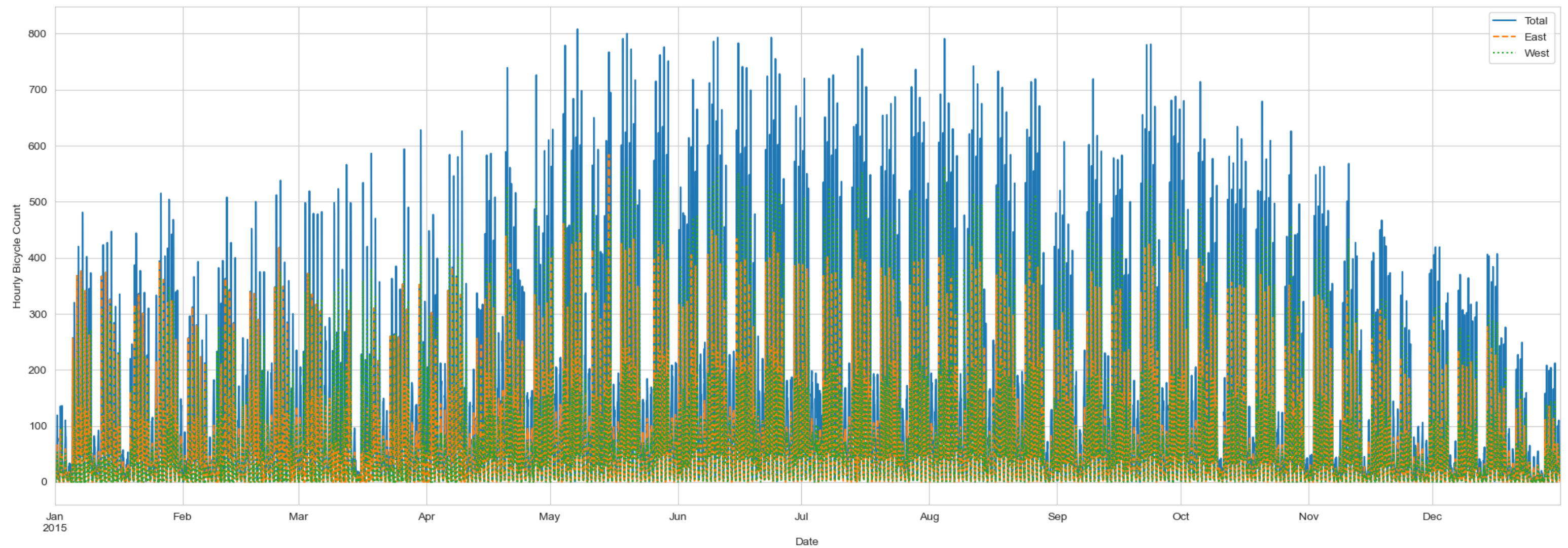
Out[65]: 'proportion missing: 0.00025'

```
In [66]:   1   df_bike_counts = df_bike_counts.fillna(method='ffill')
           2   df_bike_counts.info()
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 28440 entries, 2012-10-03 00:00:00 to 2015-12-31 23:00:00
Data columns (total 3 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   Total   28440 non-null  float64
 1   East    28440 non-null  float64
 2   West    28440 non-null  float64
dtypes: float64(3)
memory usage: 888.8 KB
```

# Plot data from 2015

# Plot data from 2015
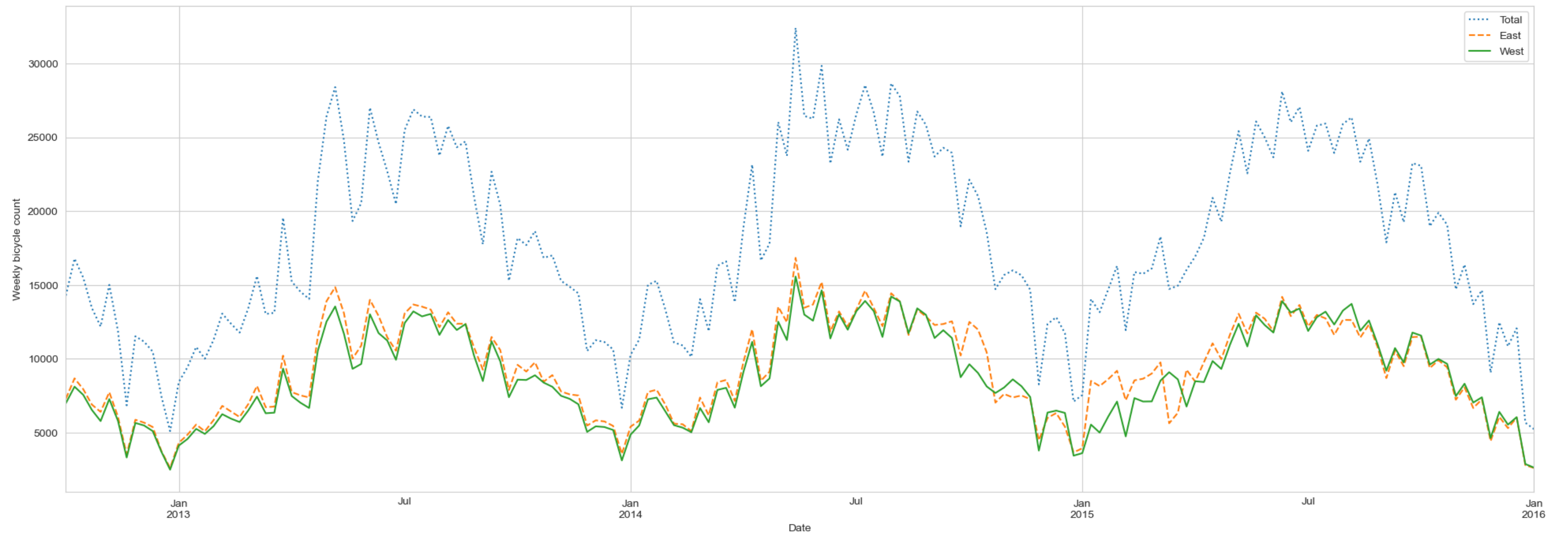
```
In [67]:  1  fig,ax = plt.subplots(1,1,figsize=(24,8))
          2  df_bike_counts.loc['2015'].plot(style=['-', '--', ':'],ax=ax)
          3  plt.ylabel('Hourly Bicycle Count');
```

# Downsample to weekly sum to smooth things out

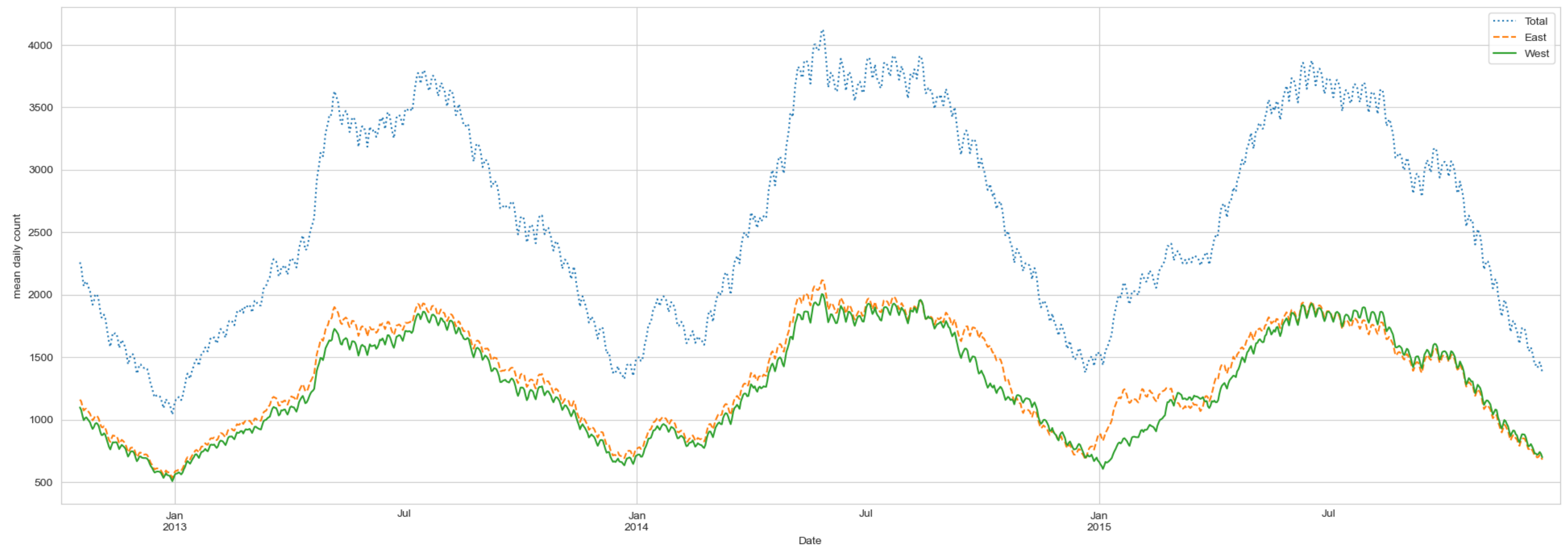# Downsample to weekly sum to smooth things out

```
In [68]:   1  weekly = df_bike_counts.resample('W').sum()
           2  weekly.plot(style=[':', '--', '-'], figsize=(24,8))
           3  plt.ylabel('Weekly bicycle count');
```

# Resample at daily for a more granular view and apply a rolling window of 30 days

# Resample at daily for a more granular view and apply a rolling window of 30 days
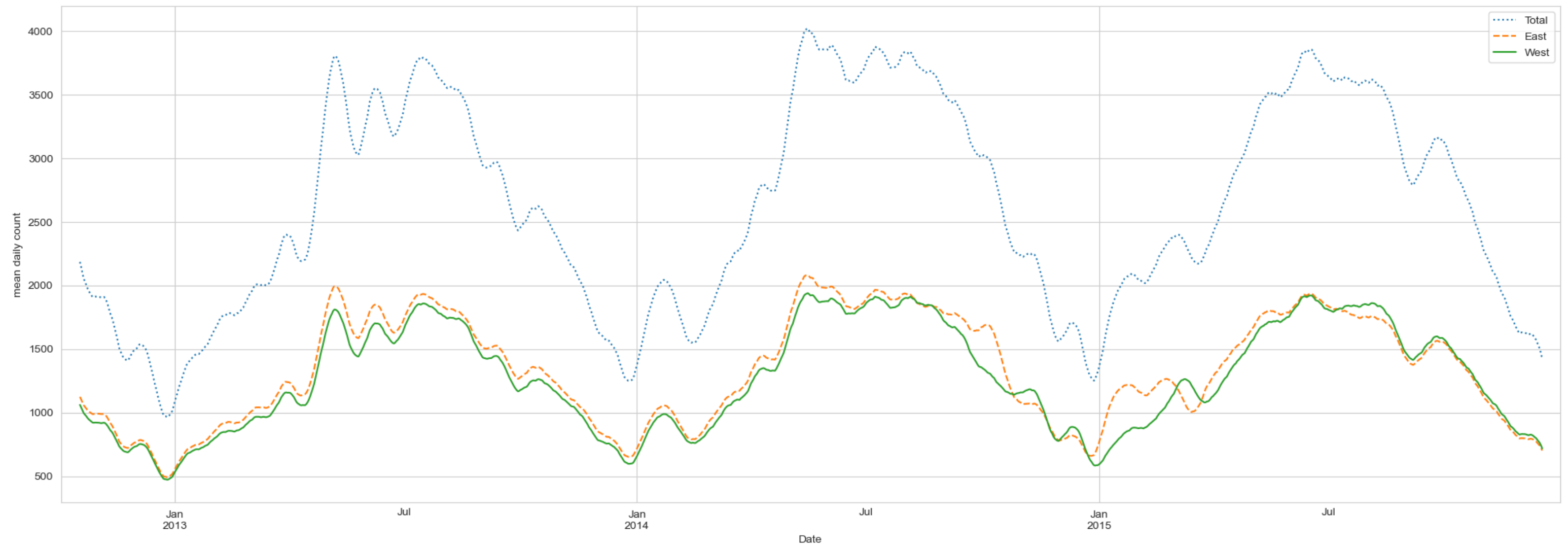
```
In [69]:  1  daily = df_bike_counts.resample('D').sum()
          2  daily.rolling(30,center=True).mean().plot(style=[':', '--', '-'], figsize=(24,8))
          3  plt.ylabel('mean daily count');
```

A wider window using a gaussian filter smooths more while accentuating daily differences

# A wider window using a gaussian filter smooths more while accentuating daily differences

```
In [70]: 1  daily.rolling(30,center=True,win_type='gaussian').mean(std=7).plot(style=[':','--','-'],figsize=(24,8));
         2  plt.ylabel('mean daily count');
```

# From Datetime to Time

# From Datetime to Time

```
In [71]:  1  #If we want to only look at time of day
          2  df_bike_counts.index.time
```

Out[71]: array([datetime.time(0, 0), datetime.time(1, 0), datetime.time(2, 0), ...,
                 datetime.time(21, 0), datetime.time(22, 0), datetime.time(23, 0)],
                dtype=object)

# From Datetime to Time

```
In [71]:  1  #If we want to only look at time of day
          2  df_bike_counts.index.time
```

```
Out[71]:  array([datetime.time(0, 0), datetime.time(1, 0), datetime.time(2, 0), ...,
                 datetime.time(21, 0), datetime.time(22, 0), datetime.time(23, 0)],
                dtype=object)
```
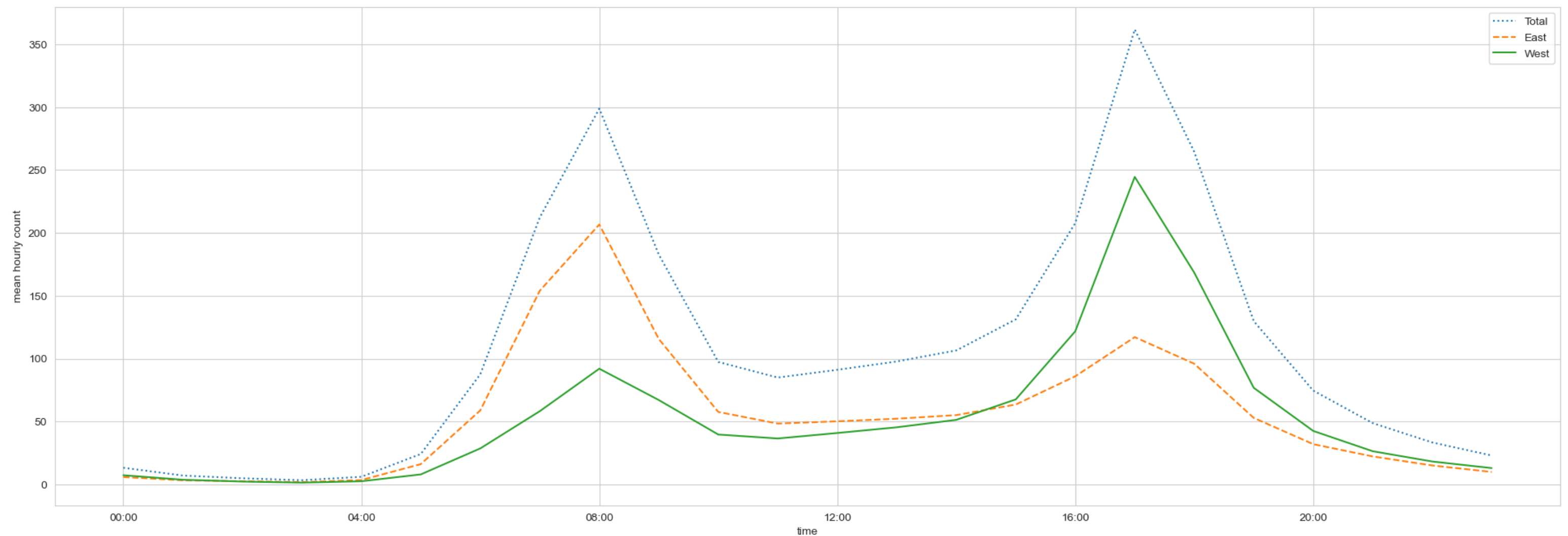
```
In [72]:  1  #Get mean data by time (hourly)
          2  by_time = df_bike_counts.groupby(df_bike_counts.index.time).mean().round(2)
          3  display(by_time)
```

|          | Total  | East   | West   |
|----------|--------|--------|--------|
| 00:00:00 | 13.34  | 5.94   | 7.40   |
| 01:00:00 | 7.15   | 3.34   | 3.81   |
| 02:00:00 | 4.97   | 2.61   | 2.36   |
| 03:00:00 | 3.43   | 1.90   | 1.52   |
| 04:00:00 | 6.13   | 3.53   | 2.59   |
| 05:00:00 | 24.26  | 16.22  | 8.04   |
| 06:00:00 | 87.65  | 58.94  | 28.71  |
| 07:00:00 | 212.38 | 154.07 | 58.31  |
| 08:00:00 | 298.85 | 206.76 | 92.09  |
| 09:00:00 | 182.88 | 115.71 | 67.17  |
| 10:00:00 | 97.45  | 57.71  | 39.74  |
| 11:00:00 | 85.06  | 48.46  | 36.61  |
| 12:00:00 | 91.21  | 50.26  | 40.96  |
| 13:00:00 | 97.83  | 52.33  | 45.50  |
| 14:00:00 | 106.61 | 55.18  | 51.42  |
| 15:00:00 | 131.29 | 63.61  | 67.67  |
| 16:00:00 | 207.88 | 86.06  | 121.82 |
| 17:00:00 | 361.78 | 117.19 | 244.59 |
| 18:00:00 | 264.50 | 96.03  | 168.47 |
| 19:00:00 | 129.85 | 52.96  | 76.89  |

# Plot by hour of the day
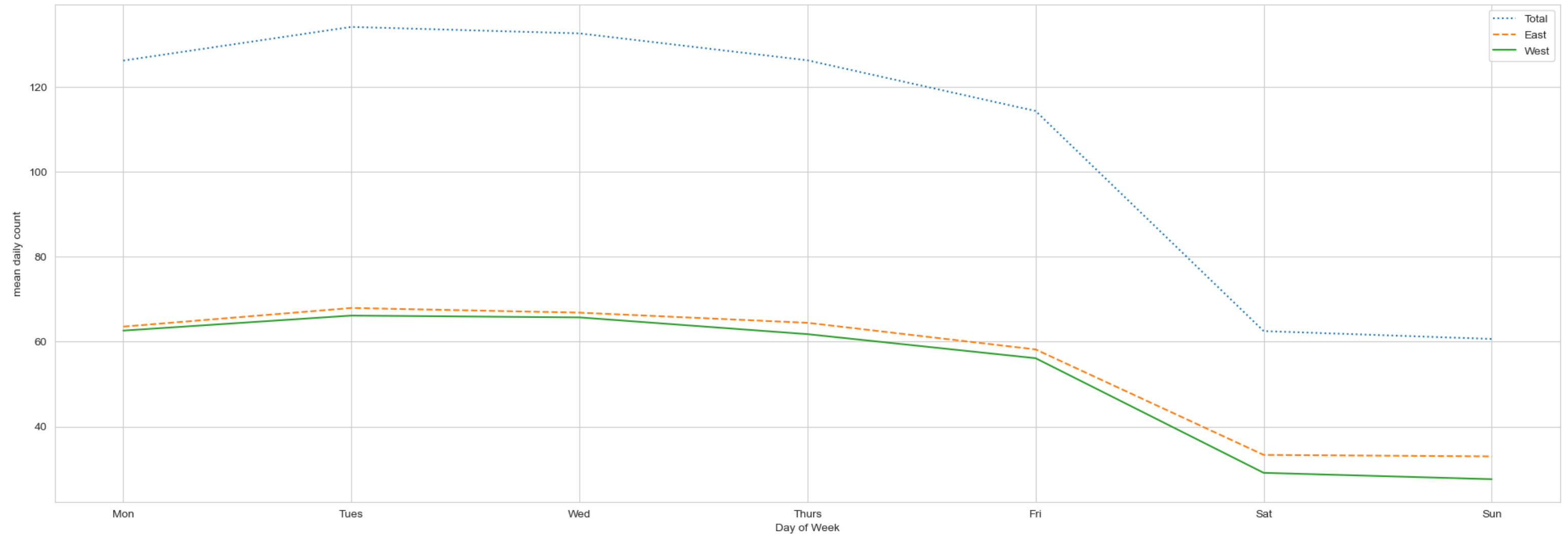
# Plot by hour of the day

```
In [73]:  1  hourly_ticks = 60 * 60 * 4 * np.arange(6)  # sec * min * every4hours
          2  by_time.plot(xticks=hourly_ticks, style=[':', '--', '-'], figsize=(24,8));
          3  plt.ylabel('mean hourly count');
```

Can also look at average by day of week
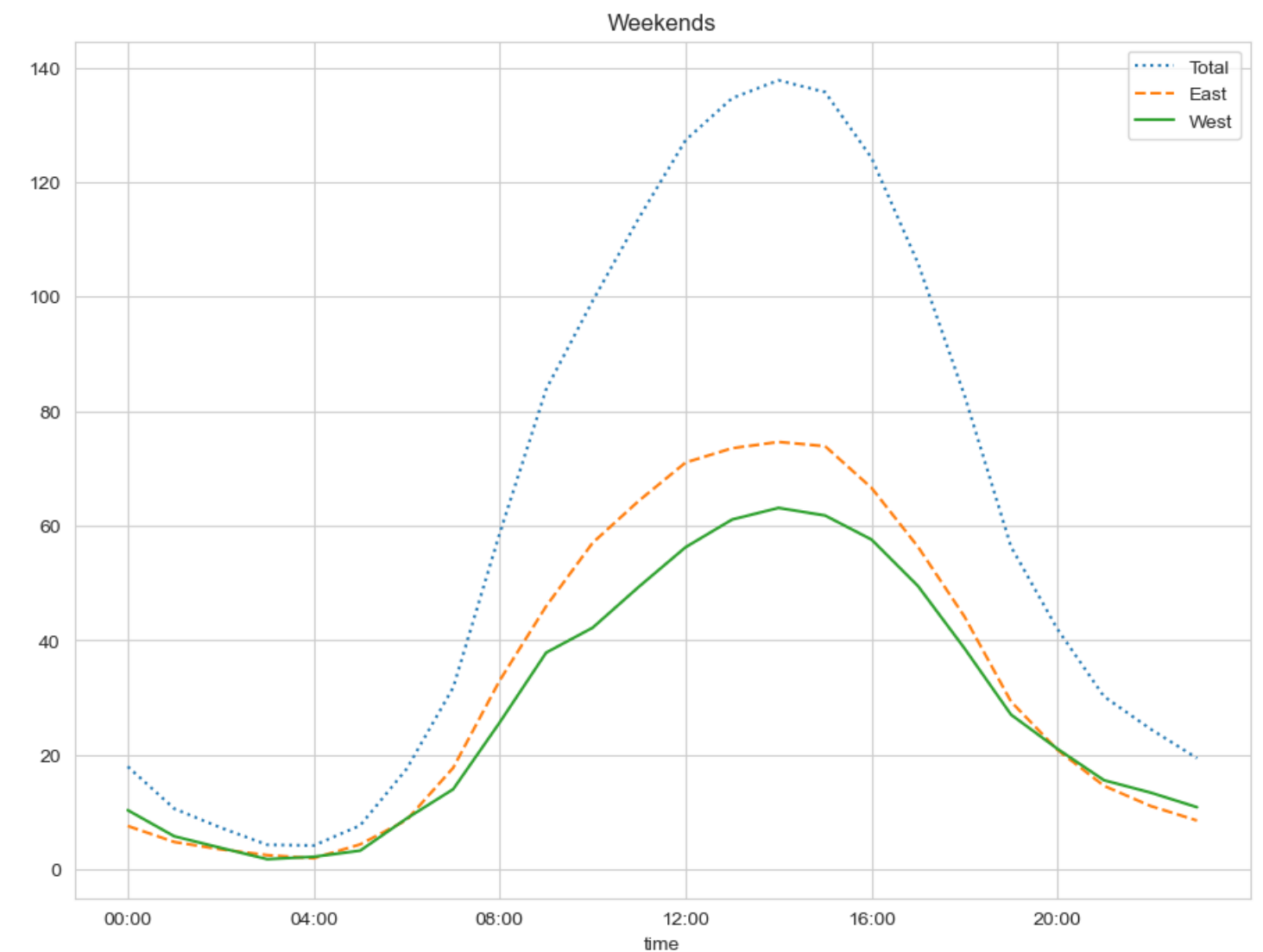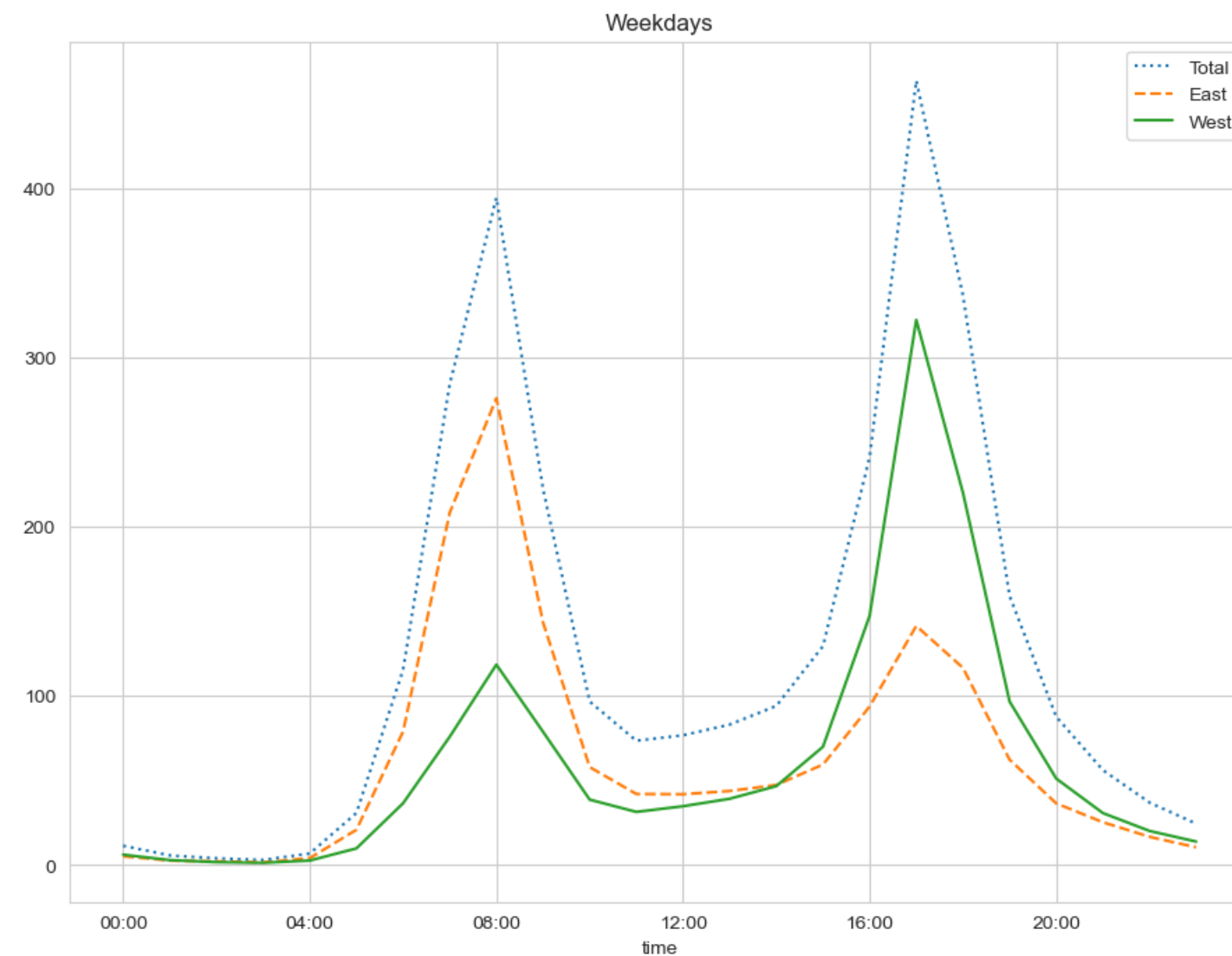
# Can also look at average by day of week

```python
# note that for dayofweek: 0 == Mon, 1 == Tues,..., 6 == 'Sun'
by_weekday = df_bike_counts.groupby(df_bike_counts.index.dayofweek).mean()
by_weekday = by_weekday.set_index(pd.Index(['Mon', 'Tues', 'Wed', 'Thurs', 'Fri', 'Sat', 'Sun']))

fig,ax = plt.subplots(1,1,figsize=(24,8))
by_weekday.plot(style=[':', '--', '-'], ax=ax);
ax.set_xlabel('Day of Week');ax.set_ylabel('mean daily count');
```

# Separate out weekdays and weekends

# Separate out weekdays and weekends

```
In [75]:  1  # create a weekend mask
          2  weekend = np.where(df_bike_counts.index.day_of_week < 5, 'Weekday', 'Weekend')
          3
          4  # get hourly mean values split by weekday, weekend
          5  by_time = df_bike_counts.groupby([weekend, df_bike_counts.index.time]).mean()
          6  fig, ax = plt.subplots(1, 2, figsize=(24, 8))
          7  by_time.loc['Weekday'].plot(ax=ax[0], title='Weekdays', xticks=hourly_ticks, style=[':', '--', '-'])
          8  by_time.loc['Weekend'].plot(ax=ax[1], title='Weekends', xticks=hourly_ticks, style=[':', '--', '-']);
```

# Can we predict daily Total bike traffic?

# Can we predict daily Total bike traffic?

```python
In [76]:  1  df_bike_counts = pd.read_csv('../data/FremontBridge_2012-2015.csv', index_col='Date', parse_dates=True)
          2  df_bike_weather = pd.read_csv('../data/BicycleWeather.csv', index_col='DATE', parse_dates=True)
          3
          4  df_bike = (
          5      df_bike_counts.loc[:,['Fremont Bridge Total']]    # keep Total as target
          6      .rename({'Fremont Bridge Total':'Total'},axis=1) # rename target column
          7      .resample('D').sum()                             # downsample to daily totals
          8  )
          9  print(df_bike.head(3))
```

```
                Total
Date
2012-10-03   3521.0
2012-10-04   3475.0
2012-10-05   3148.0
```

# Can we predict daily Total bike traffic?

```python
In [76]:  1  df_bike_counts = pd.read_csv('../data/FremontBridge_2012-2015.csv', index_col='Date', parse_dates=True)
          2  df_bike_weather = pd.read_csv('../data/BicycleWeather.csv', index_col='DATE', parse_dates=True)
          3
          4  df_bike = (
          5      df_bike_counts.loc[:,['Fremont Bridge Total']]   # keep Total as target
          6      .rename({'Fremont Bridge Total':'Total'},axis=1)  # rename target column
          7      .resample('D').sum()                              # downsample to daily totals
          8  )
          9  print(df_bike.head(3))
```

```
                Total
Date
2012-10-03   3521.0
2012-10-04   3475.0
2012-10-05   3148.0
```

# On to Feature Engineering...

# Add 'day of week'

# Add 'day of week'

```
1  day_names_map = dict(enumerate(['Mon', 'Tue', 'Wed', 'Thu', 'Fri', 'Sat', 'Sun']))
2  print(f"{day_names_map = :}")
3  df_bike['DayOfWeek'] = df_bike.index.dayofweek.map(day_names_map)
4  df_bike.head(3)
```

In [77]:

```
day_names_map = {0: 'Mon', 1: 'Tue', 2: 'Wed', 3: 'Thu', 4: 'Fri', 5: 'Sat', 6: 'Sun'}
```

Out[77]:

| Date | Total | DayOfWeek |
|---|---|---|
| 2012-10-03 | 3521.0 | Wed |
| 2012-10-04 | 3475.0 | Thu |
| 2012-10-05 | 3148.0 | Fri |

# Add 'is it a holiday' dummy feature

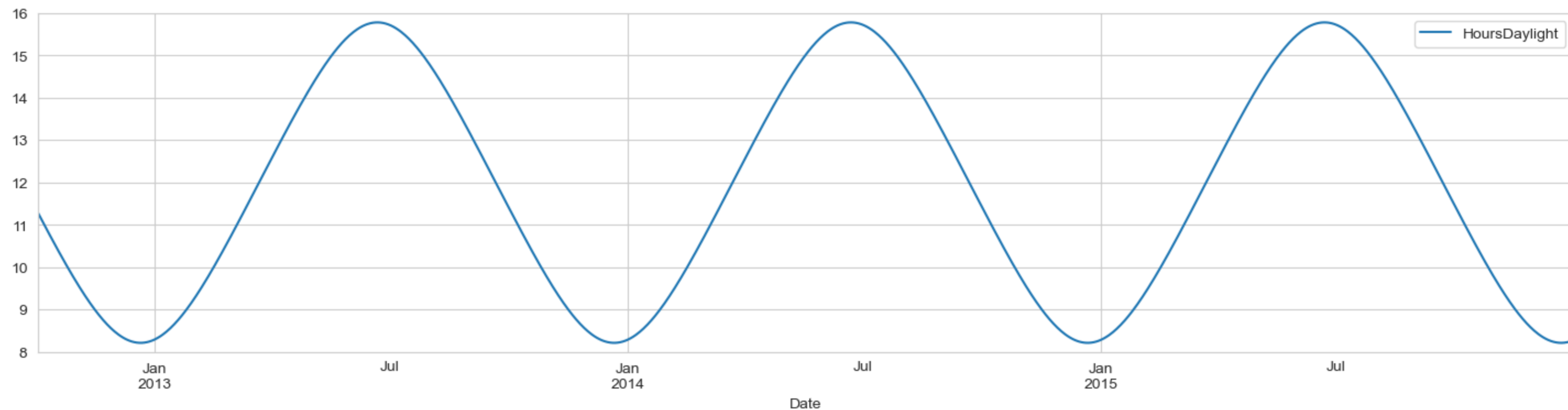# Add 'is it a holiday' dummy feature

```python
In [78]:  1  from pandas.tseries.holiday import USFederalHolidayCalendar
          2  cal = USFederalHolidayCalendar()
          3  holidays = cal.holidays('2012', '2016')
          4
          5  df_bike = df_bike.join(pd.Series(1, index=holidays, name='IsHoliday'))
          6  df_bike['IsHoliday'].fillna(0, inplace=True)
          7  print(df_bike.head(3))
```

```
                Total DayOfWeek   IsHoliday
Date
2012-10-03   3521.0      Wed          0.0
2012-10-04   3475.0      Thu          0.0
2012-10-05   3148.0      Fri          0.0
```

# Add number of hours of daylight

# Add number of hours of daylight

```python
from datetime import datetime

def hours_of_daylight(date, axis=23.44, latitude=47.61):
    """Compute the hours of daylight for the given date"""
    days = (date - datetime(2000, 12, 21)).days # days till winter solstice
    m = (1. - np.tan(np.radians(latitude))
            * np.tan(np.radians(axis) * np.cos(days * 2 * np.pi / 365.25)))
    return 24. * np.degrees(np.arccos(1 - np.clip(m, 0, 2))) / 180.

df_bike['HoursDaylight'] = list(map(hours_of_daylight, df_bike.index));

ax = df_bike[['HoursDaylight']].plot(figsize=(18,4));
ax.set_ylim(8, 16);
```

Add weather information (Q: can we predict this for future dates?)

# Add weather information (Q: can we predict this for future dates?)

In [80]:
```python
# temperatures are in 1/10 deg C; convert to C
df_bike_weather['TMIN'] /= 10
df_bike_weather['TMAX'] /= 10
df_bike_weather['TempC'] = 0.5 * (df_bike_weather['TMIN'] + df_bike_weather['TMAX'])

# precip is in 1/10 mm; convert to inches
df_bike_weather['PRCP'] /= 254
df_bike_weather['IsDryDay'] = (df_bike_weather['PRCP'] == 0).astype(int)

df_bike = df_bike.join(df_bike_weather[['PRCP', 'TempC', 'IsDryDay']],how='inner')
df_bike.head(3).round(2)
```

Out[80]:

|  | Total | DayOfWeek | IsHoliday | HoursDaylight | PRCP | TempC | IsDryDay |
|---|---|---|---|---|---|---|---|
| **2012-10-03** | 3521.0 | Wed | 0.0 | 11.28 | 0.0 | 13.35 | 1 |
| **2012-10-04** | 3475.0 | Thu | 0.0 | 11.22 | 0.0 | 13.60 | 1 |
| **2012-10-05** | 3148.0 | Fri | 0.0 | 11.16 | 0.0 | 15.30 | 1 |

# Add time of year

# Add time of year

```
In [81]: 1 df_bike['TimeOfYear'] = (df_bike.index - df_bike.index[0]).days / 365.0 # Days since the beginning of the year
         2 df_bike.head(3)
```

Out[81]:

|  | Total | DayOfWeek | IsHoliday | HoursDaylight | PRCP | TempC | IsDryDay | TimeOfYear |
|---|---|---|---|---|---|---|---|---|
| 2012-10-03 | 3521.0 | Wed | 0.0 | 11.277359 | 0.0 | 13.35 | 1 | 0.000000 |
| 2012-10-04 | 3475.0 | Thu | 0.0 | 11.219142 | 0.0 | 13.60 | 1 | 0.002740 |
| 2012-10-05 | 3148.0 | Fri | 0.0 | 11.161038 | 0.0 | 15.30 | 1 | 0.005479 |

# Generate and evaluate a model

# Generate and evaluate a model

```
In [82]:    1  from sklearn.ensemble import GradientBoostingRegressor
            2  from sklearn.dummy import DummyRegressor
            3  from sklearn.metrics import mean_absolute_error
            4
            5  # drop any rows with missing data
            6  df_bike.dropna(axis=0, how='any', inplace=True)
            7
            8  X_bike = pd.get_dummies(df_bike.loc[:,df_bike.columns != 'Total'])
            9  display(X_bike.head(1).round(2))
           10  y_bike = df_bike.Total
           11
           12  X_bike_train = X_bike.loc['2012':'2014']
           13  y_bike_train = y_bike.loc['2012':'2014']
           14  X_bike_test = X_bike.loc['2015']
           15  y_bike_test = y_bike.loc['2015']
           16
           17  dummy_bike = DummyRegressor().fit(X_bike_train,y_bike_train)
           18  gb_bike = GradientBoostingRegressor().fit(X_bike_train,y_bike_train)
           19  print(f'dummy training mae    : {mean_absolute_error(y_bike_train,dummy_bike.predict(X_bike_train)).round(2)}')
           20  print(f'one-back training mae : {mean_absolute_error(y_bike_train,y_bike_train.shift(1).fillna(0)).round(2)}')
           21  print(f'gb training set mae   : {mean_absolute_error(y_bike_train,gb_bike.predict(X_bike_train)).round(2)}')
           22  print(f'gb test set R^2       : {mean_absolute_error(y_bike_test,gb_bike.predict(X_bike_test)).round(2)}')
```
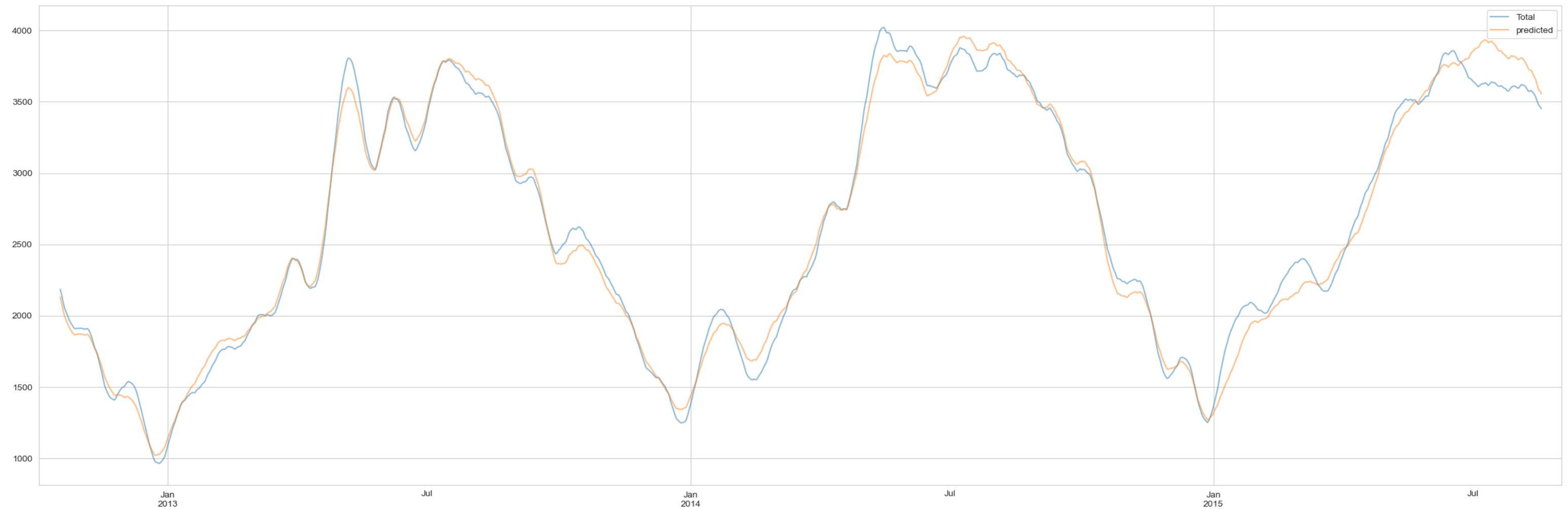
| | IsHoliday | HoursDaylight | PRCP | TempC | IsDryDay | TimeOfYear | DayOfWeek_Fri | DayOfWeek_Mon | DayOfWeek_Sat | DayOfWeek_Sun | DayOfWeek_Thu | DayOfWeek_Tue | DayOfWe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2012-10-03 | 0.0 | 11.28 | 0.0 | 13.35 | 1 | 0.0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

```
dummy training mae    : 1019.45
one-back training mae : 710.39
gb training set mae   : 213.37
gb test set R^2       : 308.51
```

# Plot predictions vs observed

# Plot predictions vs observed

```python
In [83]:  1  df_bike['predicted'] = gb_bike.predict(X_bike)
          2  df_bike[['Total', 'predicted']].rolling(30, center=True,win_type='gaussian').mean(std=7).plot(alpha=0.5,figsize=(24,8))
          3  plt.tight_layout()
```

# Time Series Operations Review

- Shifting
- Resampling
    - Downsampling
    - Upsampling
- Moving/Rolling Windows

- for more info, including time-series cross-validation:
    - sklearn: Time-related feature engineering
    - PML Chapter 13 - Modeling Sequential Data Using Recurrent Neural Network (with Tensorflow)

- for more models:
    - skforecast
    - statsmodels

# Questions re Time Series Transformations?