

Understanding the local electricity market

 The data You have access to over five years of energy price and demand data (source):

- "date" - from January 1, 2015, to October 6, 2020.
- "demand" - daily electricity demand in MWh.
- "price" - recommended retail price in AUD/MWh.
- "demand_pos_price" - total daily demand at a positive price in MWh.
- "price_positive" - average positive price, weighted by the corresponding intraday demand in AUD/MWh.
- "demand_neg_price" - total daily demand at a negative price in MWh.
- "price_negative" - average negative price, weighted by the corresponding intraday demand in AUD/MWh.
- "frac_neg_price" - the fraction of the day when the demand traded at a negative price.
- "min_temperature" - minimum temperature during the day in Celsius.
- "max_temperature" - maximum temperature during the day in Celsius.
- "solar_exposure" - total daily sunlight energy in MJ/m².
- "rainfall" - daily rainfall in mm.
- "school_day" - "Y" if that day was a school day, "N" otherwise.
- "holiday" - "Y" if the day was a state or national holiday, "N" otherwise.

Note: The price was negative during some intraday intervals, so energy producers were paying buyers rather than vice-versa.

```

1 ## 🦾 Competition challenge
2
3 Create a report that covers the following:
4
5 1. How do energy prices change throughout the year? Are there any patterns
6   by season or month of the year?
7
8 2. Build a forecast of daily energy prices the company can use as the basis
9   of its financial planning.
```

In [1]:

```
1 #conda install -c anaconda py-xgboost
```

```

In [2]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 %matplotlib inline
        6 from datetime import datetime
        7 import dateutil.parser
        8 from sklearn import linear_model
        9 from pandas.tseries.offsets import MonthEnd
       10 import warnings
       11 warnings.filterwarnings("ignore")
       12
       13 pd.set_option("expand_frame_repr", True)
       14 from sklearn.model_selection import train_test_split
       15 from sklearn.linear_model import LinearRegression
       16 from sklearn.metrics import r2_score
       17 from sklearn.metrics import mean_squared_error, mean_absolute_error
       18 from sklearn.preprocessing import MinMaxScaler
       19 import xgboost as xgb
       20
       21 col_pal = sns.color_palette()
       22 plt.style.use("fivethirtyeight")

```

```

In [3]: 1 df = pd.read_csv('energy_demand.csv', parse_dates = True)
        2 df.head()

```

```

Out[3]:

```

| | date | demand | price | demand_pos_price | price_positive | demand_neg_price | price_negati |
|---|------------|------------|-----------|------------------|----------------|------------------|--------------|
| 0 | 2015-01-01 | 99635.030 | 25.633696 | 97319.240 | 26.415953 | 2315.790 | -7.2400 |
| 1 | 2015-01-02 | 129606.010 | 33.138988 | 121082.015 | 38.837661 | 8523.995 | -47.8097 |
| 2 | 2015-01-03 | 142300.540 | 34.564855 | 142300.540 | 34.564855 | 0.000 | 0.0000 |
| 3 | 2015-01-04 | 104330.715 | 25.005560 | 104330.715 | 25.005560 | 0.000 | 0.0000 |
| 4 | 2015-01-05 | 118132.200 | 26.724176 | 118132.200 | 26.724176 | 0.000 | 0.0000 |

In [4]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2106 entries, 0 to 2105
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   date                   2106 non-null   object
1   demand                 2106 non-null   float64
2   price                  2106 non-null   float64
3   demand_pos_price       2106 non-null   float64
4   price_positive         2106 non-null   float64
5   demand_neg_price       2106 non-null   float64
6   price_negative         2106 non-null   float64
7   frac_neg_price         2106 non-null   float64
8   min_temperature        2106 non-null   float64
9   max_temperature        2106 non-null   float64
10  solar_exposure         2105 non-null   float64
11  rainfall               2103 non-null   float64
12  school_day             2106 non-null   object
13  holiday                2106 non-null   object
dtypes: float64(11), object(3)
memory usage: 230.5+ KB
```

In [5]: 1 *# A copy of the data to be used for creating a model for making the predicti*
 2 df1 = df.copy()

In [6]: 1 df['date'] = pd.to_datetime(df['date'])

In [7]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2106 entries, 0 to 2105
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   date                   2106 non-null   datetime64[ns]
1   demand                 2106 non-null   float64
2   price                  2106 non-null   float64
3   demand_pos_price       2106 non-null   float64
4   price_positive         2106 non-null   float64
5   demand_neg_price       2106 non-null   float64
6   price_negative         2106 non-null   float64
7   frac_neg_price         2106 non-null   float64
8   min_temperature        2106 non-null   float64
9   max_temperature        2106 non-null   float64
10  solar_exposure         2105 non-null   float64
11  rainfall               2103 non-null   float64
12  school_day             2106 non-null   object
13  holiday                2106 non-null   object
dtypes: datetime64[ns](1), float64(11), object(2)
memory usage: 230.5+ KB
```

In [8]: 1 df.describe()

Out[8]:

| | demand | price | demand_pos_price | price_positive | demand_neg_price | price_neg |
|-------|---------------|-------------|------------------|----------------|------------------|-------------|
| count | 2106.000000 | 2106.000000 | 2106.000000 | 2106.000000 | 2106.000000 | 2106.000000 |
| mean | 120035.476503 | 76.079554 | 119252.305055 | 76.553847 | 783.171448 | -2.611111 |
| std | 13747.993761 | 130.246805 | 14818.631319 | 130.114184 | 3578.920686 | 19.444444 |
| min | 85094.375000 | -6.076028 | 41988.240000 | 13.568986 | 0.000000 | -342.222222 |
| 25% | 109963.650000 | 38.707040 | 109246.250000 | 39.117361 | 0.000000 | 0.000000 |
| 50% | 119585.912500 | 66.596738 | 119148.082500 | 66.869058 | 0.000000 | 0.000000 |
| 75% | 130436.006250 | 95.075012 | 130119.477500 | 95.130181 | 0.000000 | 0.000000 |
| max | 170653.840000 | 4549.645105 | 170653.840000 | 4549.645105 | 57597.595000 | 0.000000 |

In [9]: 1 #checking for missing values

2
3 df.isnull().sum()

Out[9]:

| | |
|------------------|---|
| date | 0 |
| demand | 0 |
| price | 0 |
| demand_pos_price | 0 |
| price_positive | 0 |
| demand_neg_price | 0 |
| price_negative | 0 |
| frac_neg_price | 0 |
| min_temperature | 0 |
| max_temperature | 0 |
| solar_exposure | 1 |
| rainfall | 3 |
| school_day | 0 |
| holiday | 0 |
| dtype: int64 | |

In [10]: 1 #filling the missing data using mode
2 df['solar_exposure'] = df['solar_exposure'].fillna(df['solar_exposure'].mode()[0])
3 df['rainfall'] = df['rainfall'].fillna(df['rainfall'].mode()[0])

```
In [11]: 1 # Checking for duplicates
         2 df.isnull().sum()
```

```
Out[11]: date                0
         demand              0
         price               0
         demand_pos_price    0
         price_positive       0
         demand_neg_price    0
         price_negative       0
         frac_neg_price       0
         min_temperature     0
         max_temperature     0
         solar_exposure      0
         rainfall            0
         school_day          0
         holiday             0
         dtype: int64
```

```
In [12]: 1 df.duplicated().sum()
```

```
Out[12]: 0
```

```
In [13]: 1 # Splitting the date time for plotting purposes
         2
         3 df['date'] = pd.to_datetime(df['date'])
         4 df['year'] = df['date'].dt.year
         5 df['month'] = df['date'].dt.month
         6 df['day'] = df['date'].dt.day
```

```
In [14]: 1 df.columns
```

```
Out[14]: Index(['date', 'demand', 'price', 'demand_pos_price', 'price_positive',
               'demand_neg_price', 'price_negative', 'frac_neg_price',
               'min_temperature', 'max_temperature', 'solar_exposure', 'rainfall',
               'school_day', 'holiday', 'year', 'month', 'day'],
              dtype='object')
```

```
In [15]: 1 # Drop date column
         2 df.drop('date', axis = 1, inplace = True)
```

```
In [16]: 1 df.sample(3)
```

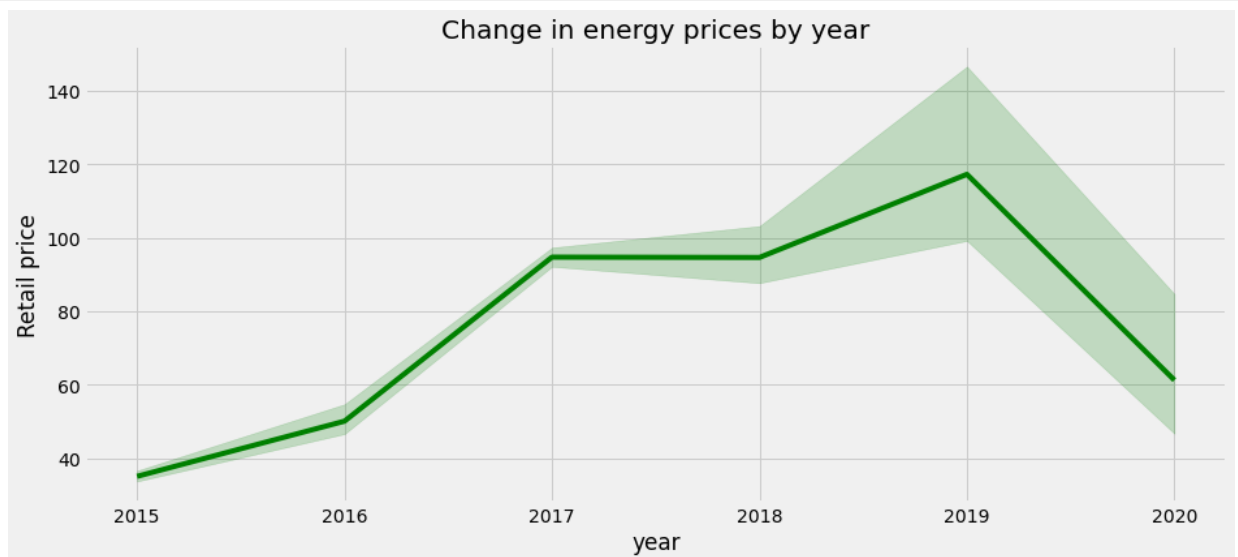
```
Out[16]:
```

| | demand | price | demand_pos_price | price_positive | demand_neg_price | price_negative |
|-------------|------------|------------|------------------|----------------|------------------|----------------|
| 1850 | 102445.880 | 50.132823 | 102445.88 | 50.132823 | 0.000 | 0.00000 |
| 1921 | 105098.025 | 23.780908 | 79509.45 | 36.334744 | 25588.575 | -15.22668 |
| 1385 | 112592.260 | 122.200670 | 112592.26 | 122.200670 | 0.000 | 0.00000 |

```
In [17]: 1 # Checking the cardinality of the year column
        2 len(df.year.unique())
```

Out[17]: 6

```
In [18]: 1 # Y axis is price closing
        2 plt.figure(figsize=(14,6))
        3 sns.lineplot(y='price', x='year', data = df, color = "green")
        4 plt.title("Change in energy prices by year")
        5 plt.ylabel("Retail price")
        6 plt.show()
```



```
In [19]: 1 my = df.groupby(['year', 'month', 'day'], as_index=False)['price'].mean()
        2 my.sample(3)
```

Out[19]:

| | year | month | day | price |
|------|------|-------|-----|-----------|
| 1226 | 2018 | 5 | 11 | 63.673070 |
| 1877 | 2020 | 2 | 21 | 48.658203 |
| 458 | 2016 | 4 | 3 | 15.804000 |

```
In [20]: 1 my2 = df.groupby(['year', 'month'], as_index=False)['price'].mean()
        2 my2.sample(3)
```

Out[20]:

| | year | month | price |
|----|------|-------|-----------|
| 15 | 2016 | 4 | 47.962939 |
| 33 | 2017 | 10 | 74.471315 |
| 21 | 2016 | 10 | 34.440483 |

```
In [21]: 1 my3 = df.groupby('month')[['price']].mean()
         2 my3.sample(3)
```

```
Out[21]:
```

| | price |
|-------|-----------|
| month | |
| 12 | 66.229981 |
| 9 | 68.578842 |
| 6 | 81.245597 |

```
In [22]: 1 my4 = df.groupby('year')[['price']].mean()
         2 my4.sample(3)
```

```
Out[22]:
```

| | price |
|------|------------|
| year | |
| 2019 | 117.281370 |
| 2017 | 94.740161 |
| 2016 | 50.094252 |

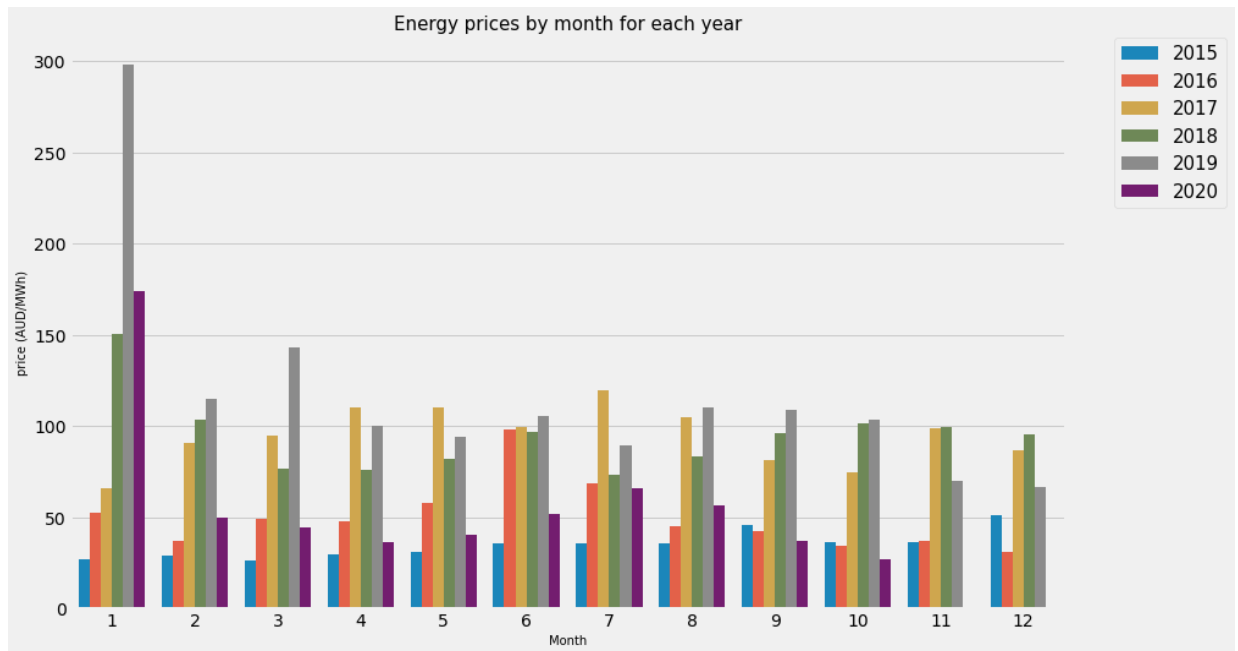
```
In [23]: 1 #add a new date time column to my2 dataframe with the date as the end of the
         2 my2['datetime'] = pd.to_datetime(my2.year.astype(str) + my2.month.astype(str)
         3 my2.sample(3))
```

```
Out[23]:
```

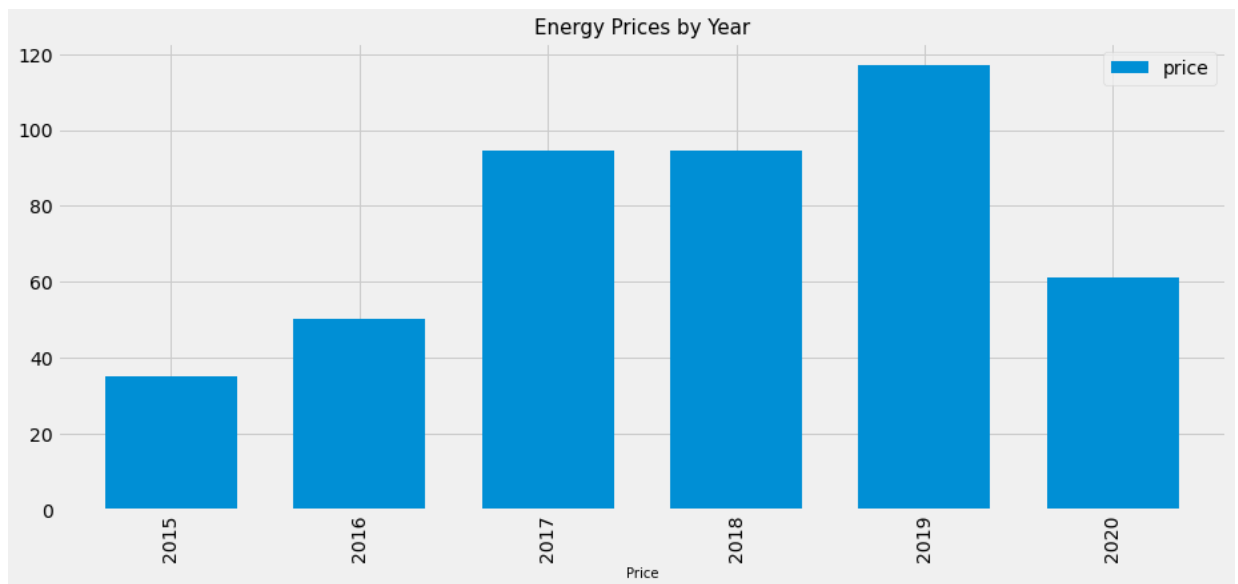
| | year | month | price | datetime |
|----|------|-------|------------|------------|
| 56 | 2019 | 9 | 109.025825 | 2019-09-30 |
| 15 | 2016 | 4 | 47.962939 | 2016-04-30 |
| 48 | 2019 | 1 | 298.171896 | 2019-01-31 |

In [24]:

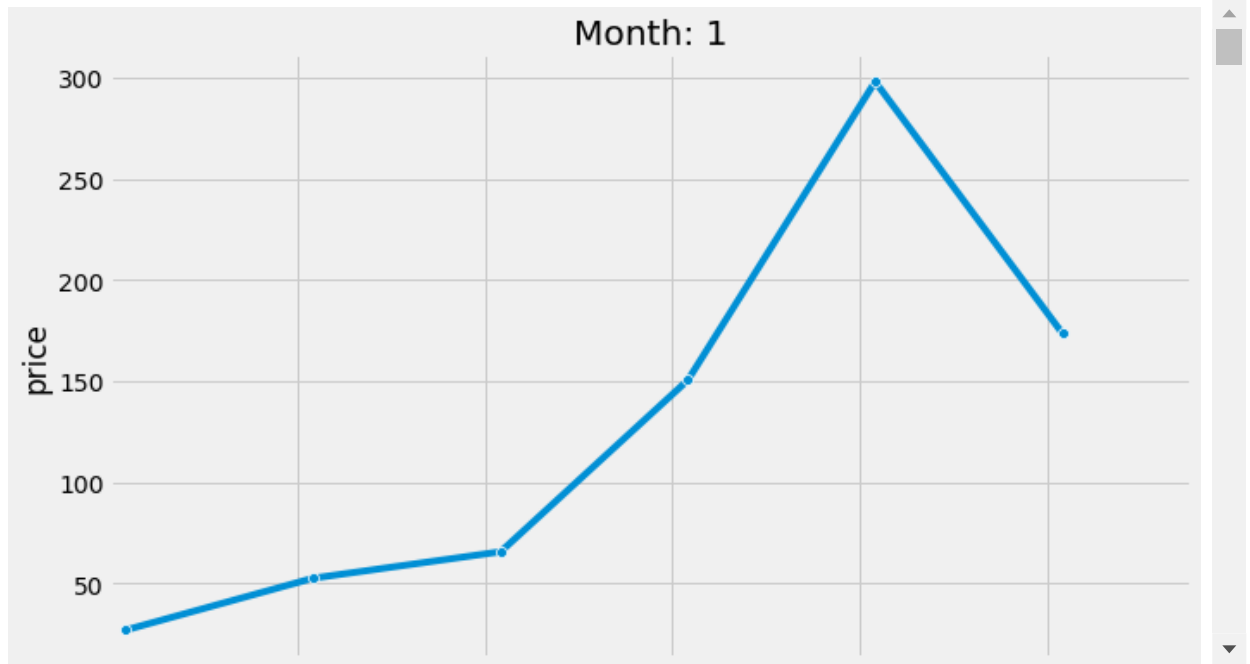
```
1 # Prices by month and year
2 plt.figure(figsize = (13, 8))
3 sns.barplot(data = my2, x = 'month', y = 'price', hue = 'year')
4 plt.legend(bbox_to_anchor = (1.05, 1), borderaxespad = 0, fontsize = 15)
5 plt.title('Energy prices by month for each year', fontsize = 15)
6 plt.ylabel('price (AUD/MWh)', fontsize = 10)
7 plt.xlabel("Month", fontsize = 10)
8 plt.show()
```




```
In [25]: 1 #Prices by year
2 my4.plot(kind = 'bar', width = 0.7, figsize = (14,6))
3 plt.title("Energy Prices by Year", fontsize = 15)
4 plt.xlabel("Year",fontsize = 10)
5 plt.xlabel("Price",fontsize = 10)
6 plt.show()
```



```
In [26]: 1 for month in my2.month.unique():
2         data = my2[my2.month == month]
3         plt.figure(figsize=(10,6))
4         sns.lineplot(data.datetime, data.price, marker = 'o')
5         plt.xlim(datetime(2015, 1, 1), datetime(2020, 10, 6))
6         plt.title(f'Month: {month}')
7         plt.ylabel('price')
8         plt.xlabel('Year')
```



seasons in australia

- #Summer - three hottest month which falls in December to February
- #Autumn - the transition month which falls in March to May
- #Winter - the three coldest months falls in June to August
- #spring - the three transition months which is September to November

Lets represent the months in numbers

- 1 = Summer (months id 12th - 2nd)
- 2 = Autumn (months id 3th - 5th)
- 3 = Winter (months id 6th - 8th)
- 4 = Spring (months id 9th - 11th)

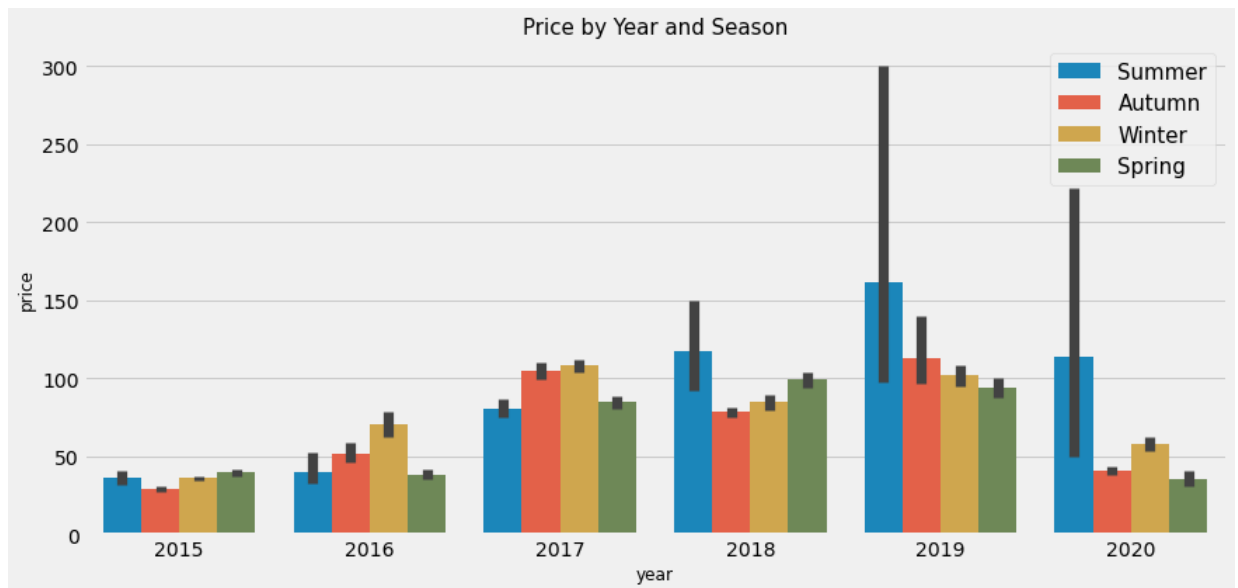
```
In [27]: 1 seasons = [1, 1, 2, 2, 2, 3, 3, 3, 4, 4, 4, 1]
2 month_to_season= dict(zip(range(1, 13), seasons))
3 df['season_id'] = df.month.map(month_to_season)
4 df['season'] = df['season_id'].map({1:'Summer', 2:'Autumn', 3:'Winter', 4:'S
5 df['rain'] = df['rainfall'] > 1
```

```
In [28]: 1 # Patterns by seasons
2 price_by_season = df.groupby(['season'], as_index = False)['price'].mean()
3 price_by_season
```

```
Out[28]:
```

| | season | price |
|---|--------|-----------|
| 0 | Autumn | 69.517789 |
| 1 | Spring | 68.451242 |
| 2 | Summer | 90.115702 |
| 3 | Winter | 76.433038 |

```
In [29]: 1 # comparing price by season..
2 plt.figure(figsize=(13,6))
3 sns.barplot(data = df, x = 'year', y = 'price', hue = 'season')
4 plt.legend(bbox_to_anchor = (1,1), fontsize = 15)
5 plt.xlabel('year', fontsize = 12)
6 plt.ylabel('price', fontsize = 12)
7 plt.title('Price by Year and Season', fontsize = 15)
8 plt.show()
```



Forecast of daily energy prices the company can use as the basis of its financial planning.

In [30]:

1 df1.head()

Out[30]:

| | date | demand | price | demand_pos_price | price_positive | demand_neg_price | price_negative |
|---|------------|------------|-----------|------------------|----------------|------------------|----------------|
| 0 | 2015-01-01 | 99635.030 | 25.633696 | 97319.240 | 26.415953 | 2315.790 | -7.2400 |
| 1 | 2015-01-02 | 129606.010 | 33.138988 | 121082.015 | 38.837661 | 8523.995 | -47.8097 |
| 2 | 2015-01-03 | 142300.540 | 34.564855 | 142300.540 | 34.564855 | 0.000 | 0.0000 |
| 3 | 2015-01-04 | 104330.715 | 25.005560 | 104330.715 | 25.005560 | 0.000 | 0.0000 |
| 4 | 2015-01-05 | 118132.200 | 26.724176 | 118132.200 | 26.724176 | 0.000 | 0.0000 |

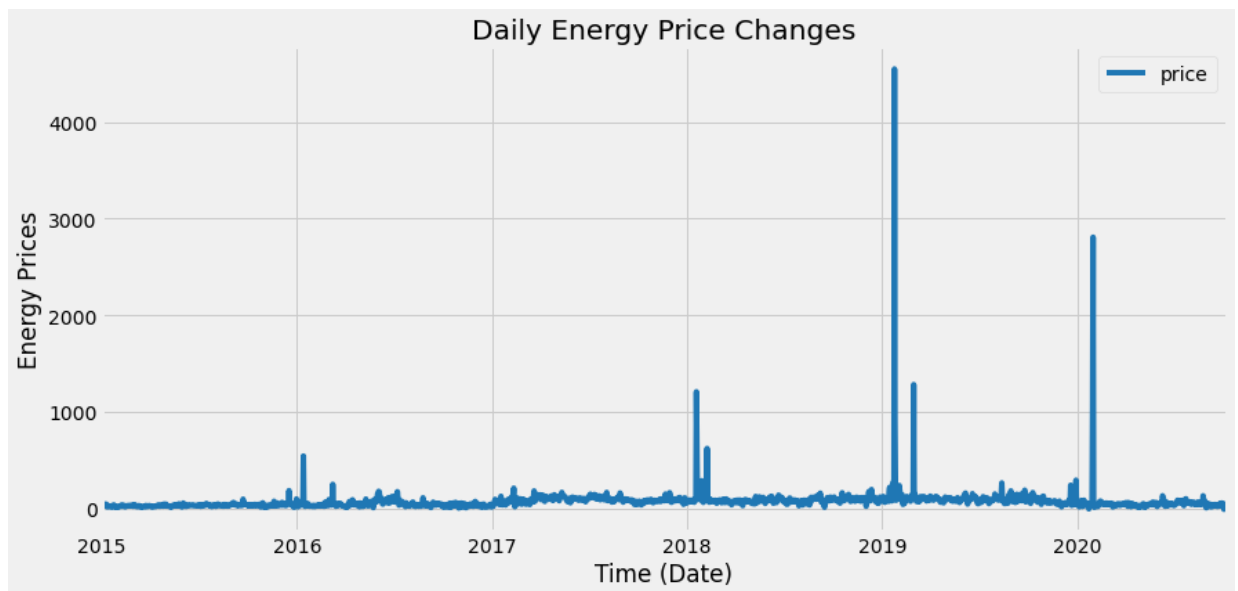
In [31]:

```
1 df1.index = pd.to_datetime(df1['date'])
2 df1.drop("date",axis = 1, inplace = True)
3 df1.sample(3)
```

Out[31]:

| | demand | price | demand_pos_price | price_positive | demand_neg_price | price_negative |
|------------|------------|-----------|------------------|----------------|------------------|----------------|
| 2016-02-19 | 122991.660 | 43.694897 | 122991.660 | 43.694897 | 0.0 | 0.0 |
| 2015-11-20 | 123412.685 | 38.026498 | 123412.685 | 38.026498 | 0.0 | 0.0 |
| 2017-11-10 | 115156.445 | 92.615537 | 115156.445 | 92.615537 | 0.0 | 0.0 |

```
In [32]: 1 # Daily energy price chnages
2 df1.plot(y = 'price', color= col_pal[0],figsize = (13,6))
3 plt.title("Daily Energy Price Changes")
4 plt.ylabel("Energy Prices")
5 plt.xlabel("Time (Date)")
6 plt.show()
```

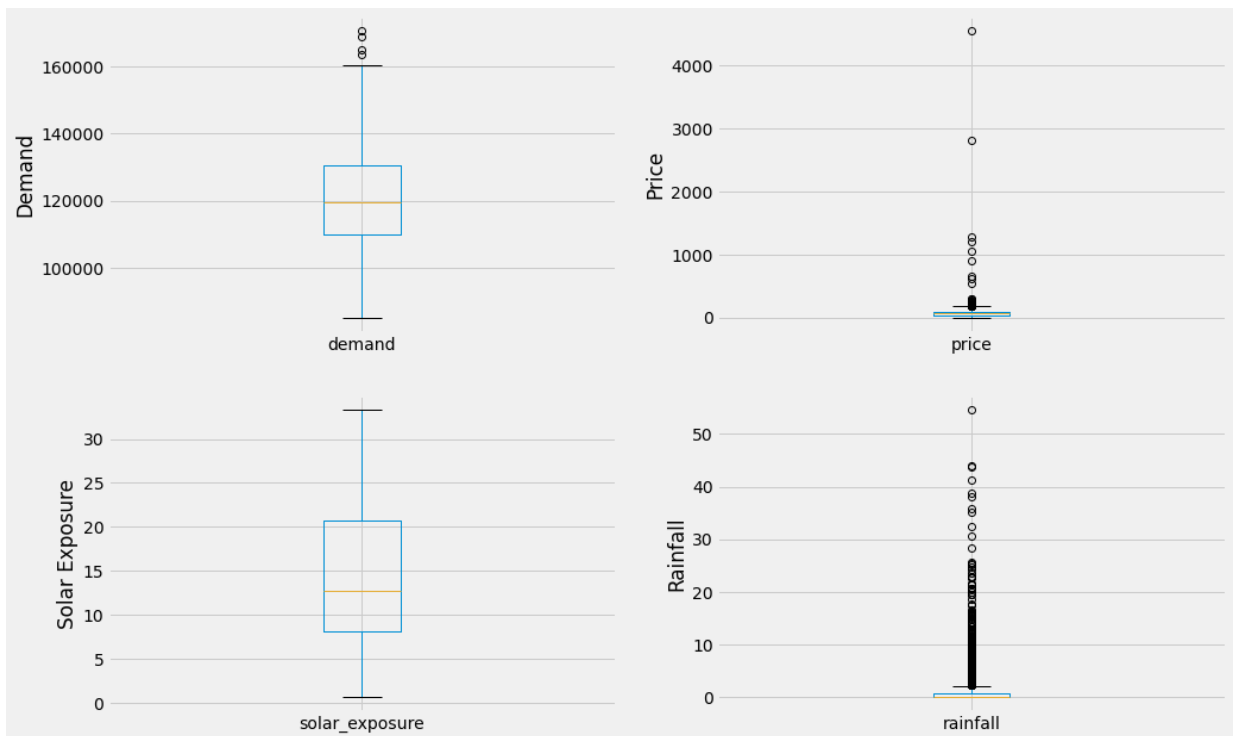


```

In [33]: 1 #Check for outliers
2 plt.figure(figsize=(15,10))
3
4 plt.subplot(2,2,1)
5 fig = df1.boxplot(column='demand')
6 fig.set_ylabel('Demand')
7
8 plt.subplot(2,2,2)
9 fig = df1.boxplot(column='price')
10 fig.set_ylabel('Price')
11
12 plt.subplot(2,2,3)
13 fig = df1.boxplot(column='solar_exposure')
14 fig.set_ylabel('Solar Exposure')
15
16 plt.subplot(2,2,4)
17 fig = df1.boxplot(column='rainfall')
18 fig.set_ylabel('Rainfall')

```

Out[33]: Text(0, 0.5, 'Rainfall')



Train and Test Split

```

In [34]: 1 len(df1.index)

```

Out[34]: 2106

```
In [35]: 1 # An estimate of around where the 70% of the data lies below(train set)
        2 df1.iloc[1475]
```

```
Out[35]: demand          152496.245
price          222.438419
demand_pos_price 152496.245
price_positive   222.438419
demand_neg_price    0.0
price_negative    0.0
frac_neg_price     0.0
min_temperature   19.4
max_temperature   30.4
solar_exposure    22.3
rainfall          0.0
school_day        N
holiday           N
Name: 2019-01-15 00:00:00, dtype: object
```

```
In [36]: 1 # splitting the date
        2 df1['year'] = df1.index.year
        3 df1['month'] = df1.index.month
        4 df1['day'] = df1.index.day
        5 df1['day of week'] = df1.index.dayofweek
```

```
In [37]: 1 # Splitting the date into test and train data
        2 split_date = '2019-01-01'
        3 train = df1[df1.index < pd.to_datetime(split_date)]
        4 test = df1[df1.index >= pd.to_datetime(split_date)]
```

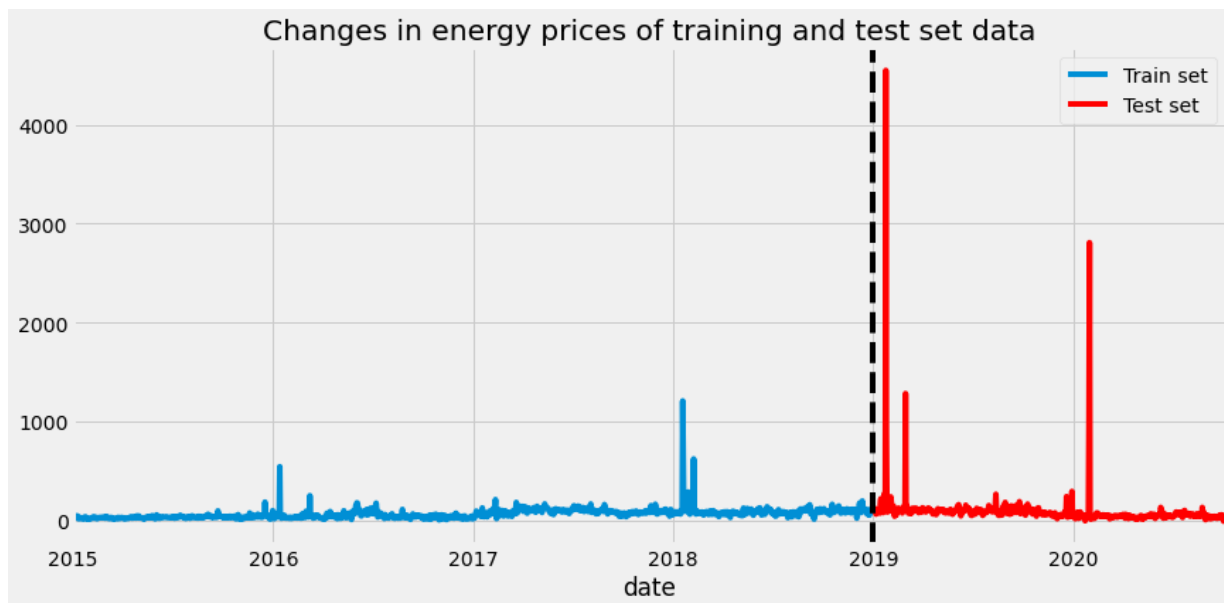
```
In [38]: 1 train.shape
```

```
Out[38]: (1461, 17)
```

```
In [39]: 1 test.shape
```

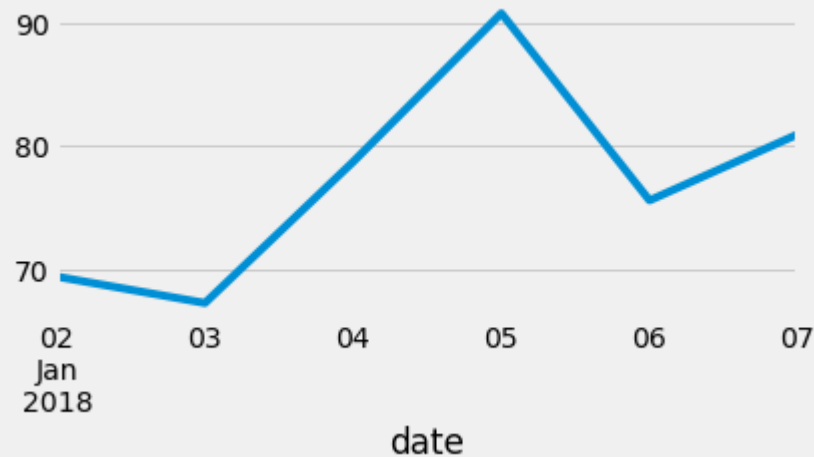
```
Out[39]: (645, 17)
```

```
In [40]: 1 fig, ax = plt.subplots(figsize=(13,6))
2 train.price.plot()
3 test.price.plot(color = "r")
4 plt.title("Changes in energy prices of training and test set data")
5 ax.axvline('2019-1-1', color = "black", ls = "--")
6 ax.legend(["Train set", "Test set"])
7 plt.show()
```

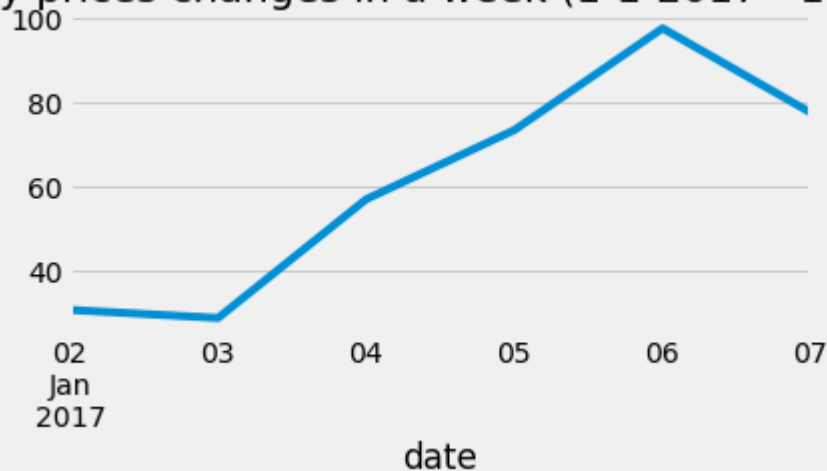


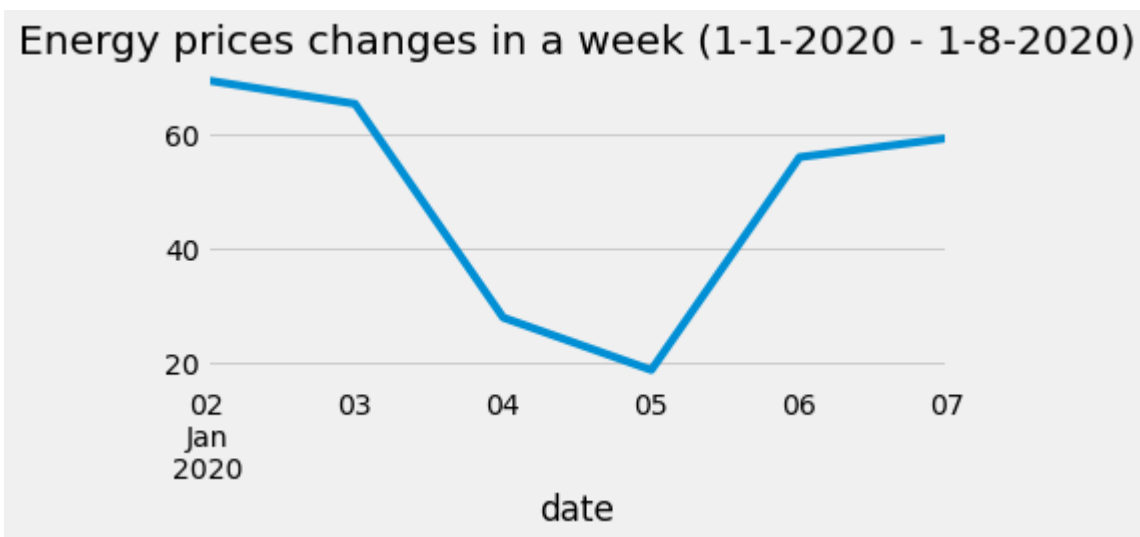

```
In [41]: 1 # Energy price changes in a random week
2 plt.figure(figsize=(10,6))
3
4 plt.subplot(2,2,1)
5 df1.loc[(df1.index > "1-1-2018") & (df1.index < "1-8-2018")].price.plot(figsize
6 plt.title("Energy prices changes in a week (1-1-2018 - 1-8-2018)")
7 plt.show()
8
9 plt.subplot(2,2,2)
10 df1.loc[(df1.index > "1-1-2017") & (df1.index < "1-8-2017")].price.plot(figsize
11 plt.title("Energy prices changes in a week (1-1-2017 - 1-8-2017)")
12 plt.show()
13
14 plt.subplot(2,2,3)
15 df1.loc[(df1.index > "1-1-2020") & (df1.index < "1-8-2020")].price.plot(figsize
16 plt.title("Energy prices changes in a week (1-1-2020 - 1-8-2020)");
```

Energy prices changes in a week (1-1-2018 - 1-8-2018)



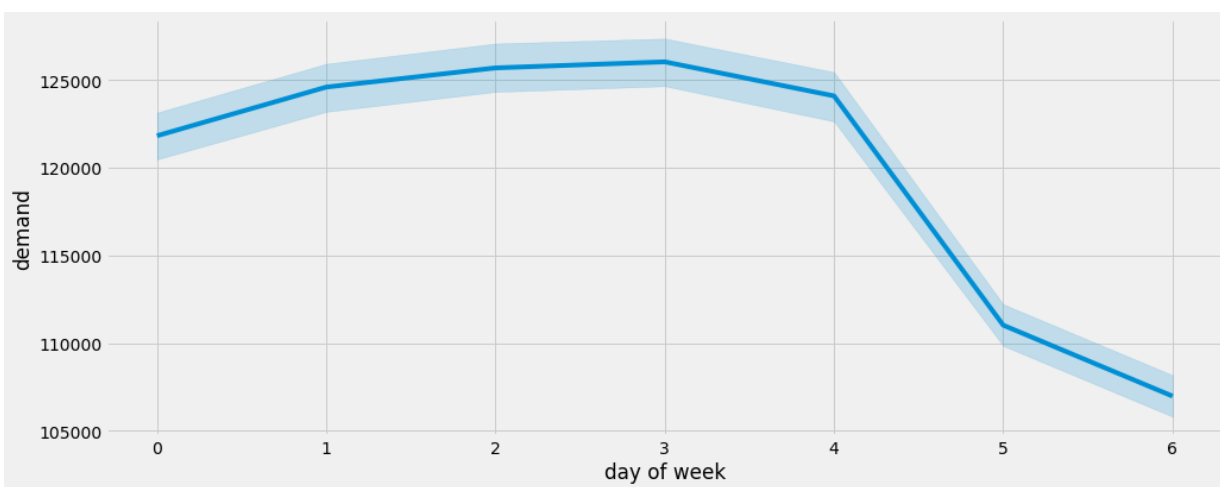
Energy prices changes in a week (1-1-2017 - 1-8-2017)





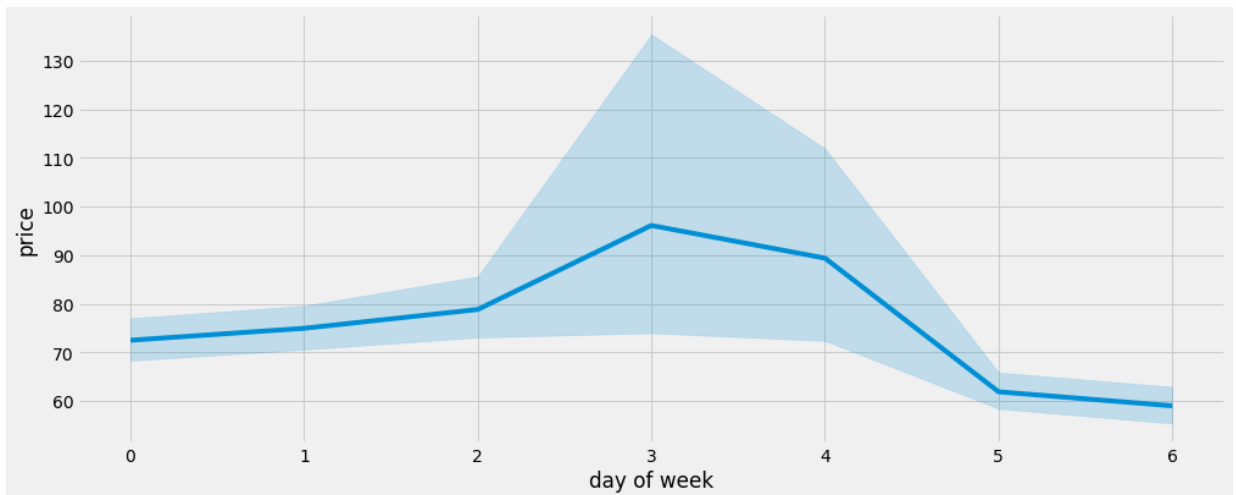
```
In [42]: 1 # Energy demand in a week
          2 plt.subplots(figsize=(15,6))
          3 sns.lineplot(data = df1, x = "day of week", y ="demand")
```

Out[42]: <AxesSubplot:xlabel='day of week', ylabel='demand'>



```
In [43]: 1 # Energy demand in a week
2 plt.subplots(figsize=(15,6))
3 sns.lineplot(data = df1, x = "day of week", y ="price")
```

Out[43]: <AxesSubplot:xlabel='day of week', ylabel='price'>



```
In [44]: 1 features = df1.drop("price", axis = 1).columns
2 print("Features: ", features)
3 target = 'price'
4 print("Target: ",target)
```

Features: Index(['demand', 'demand_pos_price', 'price_positive', 'demand_neg_p
rice',
 'price_negative', 'frac_neg_price', 'min_temperature',
 'max_temperature', 'solar_exposure', 'rainfall', 'school_day',
 'holiday', 'year', 'month', 'day', 'day of week'],
 dtype='object')
Target: price

```
In [45]: 1 X_train = train[features]
2 y_train = train[target]
3
4 X_test = test[features]
5 y_test = test[target]
```

```
In [46]: 1 #check the shape of our split data
2
3 print(X_train.shape, X_test.shape)
4 print(y_train.shape, y_test.shape)
```

(1461, 16) (645, 16)
(1461,) (645,)

```
In [47]: 1 X_train.columns
```

Out[47]: Index(['demand', 'demand_pos_price', 'price_positive', 'demand_neg_price',
 'price_negative', 'frac_neg_price', 'min_temperature',
 'max_temperature', 'solar_exposure', 'rainfall', 'school_day',
 'holiday', 'year', 'month', 'day', 'day of week'],
 dtype='object')

```
In [48]: 1 y_train = pd.DataFrame(y_train)
          2 y_test = pd.DataFrame(y_test)
```

```
In [49]: 1 def missing_values(x):
          2     return (sum(x.isna()))
          3 print("Missing values for each split: \n")
          4
          5 print("y_train: \n",y_train.apply(missing_values).where(lambda x:x!=0).dropna()
          6 print("y_test: \n",y_test.apply(missing_values).where(lambda x:x!=0).dropna()
          7 print("X_train: \n",X_train.apply(missing_values).where(lambda x:x!=0).dropna()
          8 print("X_test: \n",X_test.apply(missing_values).where(lambda x:x!=0).dropna()
```

Missing values for each split:

y_train:
Series([], dtype: float64)

y_test:
Series([], dtype: float64)

X_train:
solar_exposure 1.0
rainfall 3.0
dtype: float64

X_test:
Series([], dtype: float64)

```
In [50]: 1 # filling missing values with mode.
          2 X_train['solar_exposure'] = X_train['solar_exposure'].fillna(X_train['solar_
          3 X_train['rainfall'] = X_train['rainfall'].fillna(X_train['rainfall'].mode()[
```

In [51]: 1 X_train.info()

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 1461 entries, 2015-01-01 to 2018-12-31
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   demand                1461 non-null   float64
 1   demand_pos_price      1461 non-null   float64
 2   price_positive        1461 non-null   float64
 3   demand_neg_price      1461 non-null   float64
 4   price_negative        1461 non-null   float64
 5   frac_neg_price        1461 non-null   float64
 6   min_temperature       1461 non-null   float64
 7   max_temperature       1461 non-null   float64
 8   solar_exposure        1461 non-null   float64
 9   rainfall              1461 non-null   float64
10   school_day            1461 non-null   object
11   holiday               1461 non-null   object
12   year                 1461 non-null   int64
13   month                1461 non-null   int64
14   day                  1461 non-null   int64
15   day of week           1461 non-null   int64
dtypes: float64(10), int64(4), object(2)
memory usage: 194.0+ KB
```

In [52]:

```
1 #extraxt the categorical features
2 categorical = [var for var in X_train.columns if X_train[var].dtype=='O']
3
4 print('There are {} categorical variables \n'.format(len(categorical)))
5
6 print('They are: ', categorical)
```

There are 2 categorical variables

They are: ['school_day', 'holiday']

In [53]:

```
1 numerical = [var for var in X_train.columns if X_train[var].dtype!='O']
2
3 print('There are {} numerical variables. \n'.format(len(numerical)))
4
5 print('They are: \n', numerical)
```

There are 14 numerical variables.

They are:

['demand', 'demand_pos_price', 'price_positive', 'demand_neg_price', 'price_negative', 'frac_neg_price', 'min_temperature', 'max_temperature', 'solar_exposure', 'rainfall', 'year', 'month', 'day', 'day of week']

In [54]:

```
1 print("Unique values in school day:", X_train.school_day.unique())
2 print("Unique values in holiday:", X_train.holiday.unique())
```

Unique values in school day: ['N' 'Y']

Unique values in holiday: ['Y' 'N']

```
In [55]: 1 #encoding categorocal data
2 import category_encoders as ce
3
4 encoder =ce.BinaryEncoder(cols=['school_day', 'holiday'])
5 X_train= encoder.fit_transform(X_train)
6 X_test= encoder.fit_transform(X_test)
```

```
In [56]: 1 X_train.columns
```

```
Out[56]: Index(['demand', 'demand_pos_price', 'price_positive', 'demand_neg_price',
               'price_negative', 'frac_neg_price', 'min_temperature',
               'max_temperature', 'solar_exposure', 'rainfall', 'school_day_0',
               'school_day_1', 'holiday_0', 'holiday_1', 'year', 'month', 'day',
               'day of week'],
              dtype='object')
```

```
In [57]: 1 X_train.sample(3)
```

```
Out[57]:
```

| | demand | demand_pos_price | price_positive | demand_neg_price | price_negative | frac_neg_p |
|------------|------------|------------------|----------------|------------------|----------------|------------|
| date | | | | | | |
| 2016-06-29 | 136071.965 | 136071.965 | 80.125243 | 0.0 | 0.0 | |
| 2018-06-13 | 131574.605 | 131574.605 | 84.579533 | 0.0 | 0.0 | |
| 2018-11-22 | 107278.830 | 107278.830 | 77.208903 | 0.0 | 0.0 | |

In [58]: 1 X_train.info()

```

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 1461 entries, 2015-01-01 to 2018-12-31
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   demand                1461 non-null   float64
 1   demand_pos_price      1461 non-null   float64
 2   price_positive        1461 non-null   float64
 3   demand_neg_price      1461 non-null   float64
 4   price_negative        1461 non-null   float64
 5   frac_neg_price        1461 non-null   float64
 6   min_temperature       1461 non-null   float64
 7   max_temperature       1461 non-null   float64
 8   solar_exposure        1461 non-null   float64
 9   rainfall              1461 non-null   float64
10   school_day_0          1461 non-null   int64   
11   school_day_1          1461 non-null   int64   
12   holiday_0             1461 non-null   int64   
13   holiday_1             1461 non-null   int64   
14   year                  1461 non-null   int64   
15   month                 1461 non-null   int64   
16   day                   1461 non-null   int64   
17   day of week           1461 non-null   int64   
dtypes: float64(10), int64(8)
memory usage: 216.9 KB

```

In [59]: 1 X_test.head()

Out[59]:

| | demand | demand_pos_price | price_positive | demand_neg_price | price_negative | frac_neg_p |
|------------|------------|------------------|----------------|------------------|----------------|------------|
| date | | | | | | |
| 2019-01-01 | 98933.060 | 98933.060 | 78.560979 | 0.0 | 0.0 | |
| 2019-01-02 | 106470.675 | 106470.675 | 92.202011 | 0.0 | 0.0 | |
| 2019-01-03 | 118789.605 | 118789.605 | 127.380303 | 0.0 | 0.0 | |
| 2019-01-04 | 133288.460 | 133288.460 | 121.020997 | 0.0 | 0.0 | |
| 2019-01-05 | 97262.790 | 97262.790 | 83.493520 | 0.0 | 0.0 | |

In [60]:

```

1 X_test.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 645 entries, 2019-01-01 to 2020-10-06
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   demand                645 non-null    float64
 1   demand_pos_price      645 non-null    float64
 2   price_positive        645 non-null    float64
 3   demand_neg_price      645 non-null    float64
 4   price_negative        645 non-null    float64
 5   frac_neg_price        645 non-null    float64
 6   min_temperature       645 non-null    float64
 7   max_temperature       645 non-null    float64
 8   solar_exposure        645 non-null    float64
 9   rainfall              645 non-null    float64
10   school_day_0          645 non-null    int64
11   school_day_1          645 non-null    int64
12   holiday_0             645 non-null    int64
13   holiday_1             645 non-null    int64
14   year                  645 non-null    int64
15   month                 645 non-null    int64
16   day                   645 non-null    int64
17   day of week           645 non-null    int64
dtypes: float64(10), int64(8)
memory usage: 95.7 KB

```

In [61]:

```

1 #feature scaling`
2 scaler= MinMaxScaler()
3
4 X_train = scaler.fit_transform(X_train)
5 X_test = scaler.fit_transform(X_test)

```



```
In [62]: 1 # Model
2 # an estimation of 1000 trees to be created
3 model = xgb.XGBRegressor(n_estimator = 1000, early_stopping_rounds = 50, lea
4 model.fit(X_train, y_train, eval_set= [(X_train, y_train), (X_test, y_test)],
```

[17:58:38] WARNING: C:/buildkite-agent/builds/buildkite-windows-cpu-autoscaling-group-i-03de431ba26204c4d-1/xgboost/xgboost-ci-windows/src/learner.cc:767: Parameters: { "n_estimator" } are not used.

| | | |
|------|----------------------------|-----------------------------|
| [0] | validation_0-rmse:78.11403 | validation_1-rmse:234.96259 |
| [10] | validation_0-rmse:31.85523 | validation_1-rmse:206.09387 |
| [20] | validation_0-rmse:14.84897 | validation_1-rmse:192.16578 |
| [30] | validation_0-rmse:7.77317 | validation_1-rmse:184.57245 |
| [40] | validation_0-rmse:4.39702 | validation_1-rmse:180.38736 |
| [50] | validation_0-rmse:2.58362 | validation_1-rmse:178.04991 |
| [60] | validation_0-rmse:1.54986 | validation_1-rmse:176.66196 |
| [70] | validation_0-rmse:0.94759 | validation_1-rmse:175.86546 |
| [80] | validation_0-rmse:0.59968 | validation_1-rmse:175.39540 |
| [90] | validation_0-rmse:0.40597 | validation_1-rmse:175.11223 |
| [99] | validation_0-rmse:0.31197 | validation_1-rmse:174.95185 |

```
Out[62]: XGBRegressor(base_score=0.5, booster='gbtree', callbacks=None,
                    colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,
                    early_stopping_rounds=50, enable_categorical=False,
                    eval_metric=None, feature_types=None, gamma=0, gpu_id=-1,
                    grow_policy='depthwise', importance_type=None,
                    interaction_constraints='', learning_rate=0.1, max_bin=256,
                    max_cat_threshold=64, max_cat_to_onehot=4, max_delta_step=0,
                    max_depth=6, max_leaves=0, min_child_weight=1, missing=nan,
                    monotone_constraints='()', n_estimator=1000, n_estimators=100,
                    n_jobs=0, num_parallel_tree=1, predictor='auto', ...)
```

Forecast on the test data

```
In [63]: 1 y_pred = model.predict(X_test)
```

```
In [64]: 1 y_pred = pd.DataFrame(y_pred, columns = ['Predicted energy prices'])
2 y_pred.head(3)
```

```
Out[64]: Predicted energy prices
```

| | |
|---|-----------|
| 0 | 30.785151 |
| 1 | 34.032856 |
| 2 | 43.501045 |

In [65]: 1 y_test.head(3)

Out[65]:

| | price |
|------------|------------|
| date | |
| 2019-01-01 | 78.560979 |
| 2019-01-02 | 92.202011 |
| 2019-01-03 | 127.380303 |

In [66]: 1 mse = mean_squared_error(y_test,y_pred)
2 rmse = np.sqrt(mse)
3 rmse

Out[66]: 174.9518464452088

In [67]: 1 mae = mean_absolute_error(y_test,y_pred)
2 mae

Out[67]: 59.073529425444185

In [68]: 1 df1.index.max()

Out[68]: Timestamp('2020-10-06 00:00:00')

forecasting prices for the next one year

In [69]: 1 df1.index.max()

Out[69]: Timestamp('2020-10-06 00:00:00')

In [70]: 1 *# Creating a future data frame of one year for each of the days*
2 dates = pd.date_range('2020-10-7', '2021-10-6', freq = '1d')
3 dates

Out[70]: DatetimeIndex(['2020-10-07', '2020-10-08', '2020-10-09', '2020-10-10',
'2020-10-11', '2020-10-12', '2020-10-13', '2020-10-14',
'2020-10-15', '2020-10-16',
...,
'2021-09-27', '2021-09-28', '2021-09-29', '2021-09-30',
'2021-10-01', '2021-10-02', '2021-10-03', '2021-10-04',
'2021-10-05', '2021-10-06'],
dtype='datetime64[ns]', length=365, freq='D')

```
In [71]: 1 dates_df = pd.DataFrame(index =dates)
          2 dates_df.head()
```

```
Out[71]:
```

```
2020-10-07
2020-10-08
2020-10-09
2020-10-10
2020-10-11
```

```
In [ ]: 1
```