

Imports

In [1]:

```
import os

import pandas as pd
pd.options.mode.chained_assignment = None

import datetime as dt
import numpy as np
import statsmodels.api as sm

import matplotlib.pyplot as plt
import matplotlib.dates as mdates
import seaborn as sns
plt.style.context(("seaborn", "ggplot"))
os.chdir('..')
```

In [2]:

```
assert os.getcwd() == 'C:\\Users\\Garrett\\OneDrive\\Projects\\TempPred'
```

Goal: Predict peak-load on Thursday June 18, 2020 and hour of peak-load

In [3]:

```
def plot_dual_tseries(data_df, ts_1, ts_2, ts_1_title, ts_2_title, plot_title):
    """
    =====
    Description:
        - Convenience function for generating consistent dual plots of time series
    -----
    Inputs:
        - data_df: Dataframe containing the two time series to be plotted
        - ts_1, ts_2: column names from data_df which denote the respective series
        - ts_1_title, ts_2_title: strings to be used as the respective y-axis labels
        - plot_title: string giving the overall plot name
    Outputs:
        - None; plot generated on call
    -----
    Notes:
        - data_df index should be set as a Datetime index prior to being passed into the function
        - Should add an assert statement here in future?
        - Question is how to handle different freq.
    =====
    """

fig1, ax1 = plt.subplots(figsize = (20, 10))

ax1.plot(ts_1, data = data_df, color = 'blue', linewidth = 2)
ax1.tick_params(axis = 'y', labelcolor = 'blue', labelsize = 15)
```

```

locator = mdates.AutoDateLocator(minticks=36)
ax1.xaxis.set_major_locator(locator)
ax1.xaxis.set_tick_params(rotation = 45, labelsize = 15)
ax1.set_xlim(data_df.index[0], data_df.index[-1])
ax1.set_ylabel(ts_1_title, fontsize = 20, color = 'blue')

ax2 = ax1.twinx()
ax2.plot(ts_2, data = data_df, color = 'red', linewidth = 2)
ax2.tick_params(axis = 'y', labelcolor = 'red', labelsize = 15)
ax2.set_ylabel(ts_2_title, fontsize = 20, color = 'red')

fig1.suptitle(plot_title, fontsize = 25)

fig1.tight_layout()

def plot_annual_loadprofile(data_df, month_num):
    """
    =====
    Description:
        - Generate scatterplot of load vs. temperature annually, colored by peak vs. off-peak periods
    -----
    Inputs:
        - data_df: Dataframe containing load and temperature data
        - month_num: integer corresponding to the month in question
    Outputs:
        - None; plot generated by call
    -----
    Notes:
        - Would refine by altering subplot titles, possibly legend as well
    =====
    """
    month_names = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sept', 'Oct', 'Nov', 'Dec']

    df_tmp = data_df.loc[df['mth'] == month_num]
    month_name = month_names[month_num-1] #Because list is 0 indexed, but month number in df starts at 1

    plt.figure(figsize = (20,10))
    sns.set_context('talk')
    g1 = sns.relplot(data = df_tmp, x = 'temp', y = 'load',
                      hue = 'on_peak', palette = 'deep', col = 'yr',
                      height = 8, aspect=.75)

    g1.set_axis_labels('Temperature (F)', "Load (MW)")
    plt.suptitle(month_name+' Load vs Temperature by Year, Load Type', fontsize = 20)
    plt.tight_layout()

def plot_monthly_loadprofile(data_df):
    """
    =====
    Description:
        - Generate scatterplot of load vs. temperature monthly, colored by peak vs. off-peak periods
    -----
    Inputs:
        - data_df: Dataframe containing load and temperature data
    Outputs:
        - None; plot generated by call
    -----
    Notes:

```

```

- Need to move legend in future
=====

"""

plt.figure(figsize = (20,10))
month_names = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sept', 'Oct', 'Nov', 'Dec']

sns.set_context('talk')
g2 = sns.relplot(data = data_df, x = 'temp', y = 'load',
                  hue = 'on_peak', palette = 'deep', col = 'mth', col_wrap = 4,
                  height = 8, aspect=.75)

for n in g2.axes_dict.keys():
    tmp = g2.axes_dict[n]
    mth = month_names[n-1]
    tmp.set_title(mth, fontsize = 20)

g2.set_axis_labels('Temperature (F)', "Load (MW)")
plt.suptitle('Load vs Temperature by Month, Load Type', fontsize = 20)
plt.tight_layout()

def plot_dayofweek_loadprofile(data_df, month_num):

"""

Description:
- Generate scatterplot and violinplots of load vs. temperature daily
-----
Inputs:
- data_df: Dataframe containing load and temperature data
- month_num: integer corresponding to the month in question
Outputs:
- None; plot generated by call
-----
Notes:
- Need to move legend in future
=====

"""

month_names = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sept', 'Oct', 'Nov', 'Dec']
day_names = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']

df_tmp = data_df.loc[df['mth'] == month_num]
month_name = month_names[month_num-1] #Because list is 0 indexed, but month number in df starts at 1

plt.figure(figsize = (20,10))
sns.set_context('talk')
g3 = sns.relplot(data = df_tmp, x = 'temp', y = 'load',
                  hue = 'on_peak', palette = 'deep', col = 'dow',
                  height = 8, aspect=.75, col_wrap = 4)

for n in range(len(g3.axes_dict.keys())):
    g3.axes_dict[n].set_title(day_names[n], fontsize = 20)

g3.set_axis_labels('Temperature (F)', "Load (MW)")
plt.suptitle(month_name+' Load vs Temperature by Day of Week, Load Type', fontsize = 20)
plt.tight_layout()

plt.figure(figsize = (20,10))
sns.set_context('talk')

```

```
g4 = sns.violinplot(data = df_tmp, x = 'dow', y = 'load', hue = 'on_peak', split = True, inner = 'quart')

plt.suptitle(month_name+' Load Distribution by Day of Week, Load Type', fontsize = 20)
g4.set_xlabel(None)
g4.set_ylabel("Load (MW)")
g4.set_xticklabels(day_names)
plt.tight_layout()

def plot_hourly_loadprofile(data_df, month_num):
    """
    =====
    Description:
        - Generate scatterplot and violinplots of load vs. temperature hourly
    -----
    Inputs:
        - data_df: Dataframe containing load and temperature data
        - month_num: integer corresponding to the month in question
    Outputs:
        - None; plot generated by call
    -----
    Notes:
        - Need to move legend in future
    =====
    """
    month_names = ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sept', 'Oct', 'Nov', 'Dec']
    df_tmp = data_df.loc[df['mth'] == month_num]
    month_name = month_names[month_num-1] #Because list is 0 indexed, but month number in df starts at 1

    plt.figure(figsize = (20,10))
    sns.set_context('talk')
    g1 = sns.violinplot(data = df_tmp, x = 'hr', y = 'load', inner = 'quart')

    plt.suptitle(month_name+' Load Distribution by Hour, Load Type', fontsize = 20)
    g1.set_xlabel(None)
    g1.set_ylabel("Load (MW)")

    plt.tight_layout()

    # If I were moving this code forward and operationalizing it, I would likely opt to create a wrapper function
    # to the effect of .plot_loadprofile(data_df, month_num, profile_type) and aggregate above functions
    # or create a class, but I think breaking out individually is helpful to better flow with the markdown discussion
```

In [4]:

```
def gen_peak_id(df, ts_col, hol_df, hol_col):
    """
    =====
    Description:
        - Function for parsing load into on-peak and off-peak, accounting for business days and that year's holida
        - On-peak is defined as hours ending 7-22 on non-NERC holiday business days, and off-peak is all other
    -----
    Inputs:
        - df: Pandas Dataframe containing the load data and dates
        - ts_col: name of the column corresponding to the date (hourly index)
        - hol_df: Dataframe containing the holiday dates for the timespan covered by the data
        - hol_col: name of the column in hol_df dataframe containing the dates
    Outputs:
        - df_tmp: copy of input Dataframe, enhanced with an indicator variable for on-peak periods
```

Notes:

- Occasionally Excel hourly data reads in strangely (23:59 instead of 00:00 for example), so we round to the hour.
- Referenced against <https://www.energygps.com/HomeTools/PowerCalendar>

.....

```
dt_id = pd.DatetimeIndex(df[ts_col].round('H'), normalize = True)
hday_idx = pd.DatetimeIndex(hol_df[hol_col], normalize=True)
```

```
df_tmp = df.copy()
df_tmp['ts'], df_tmp['yr'], df_tmp['mth'], df_tmp['dom'], df_tmp['dow'], df_tmp['hr'] = dt_id, dt_id.year, dt_id.month,
df_tmp['nerc_hol'] = np.nan
```

```
df_tmp.loc[dt_id.isin(hday_idx), 'nerc_hol'] = 1
df_tmp['nerc_hol'].ffill(limit = 23, inplace = True)
df_tmp['nerc_hol'].fillna(0, inplace = True)
```

```
peak_mask = (df_tmp['nerc_hol'] == 0) & (df_tmp['dow'] < 5) & ((df_tmp['hr'] > 6) & (df_tmp['hr']<23))
```

```
df_tmp.loc[peak_mask, 'on_peak'] = 1
df_tmp.loc[~peak_mask, 'on_peak'] = 0
```

```
return df_tmp.drop(columns = 'nerc_hol')
```

In [5]:

```
class LoadModel:
    def __init__(self, data_df, y_col, x_cols):
        ....
```

Description:

- This class is used to specify, fit, and evaluate the various alternative regression models utilized in this notebook.
- The overall class consists of 2 function calls that handle most of the work:
 - LoadModel.fit_model(): allows user to specify type of regression (OLS or Robust) as well as whether heteroscedasticity aware standard errors should be used.
 - LoadModel.plot_modelfit(): leverages hourly hues to plot the residuals of the regression vs. temperature.

Inputs:

- data_df: Dataframe containing the data used to build the model
- y_col: string indicating what column to use as a dependent variable
- x_cols: list or string indicating what column(s) to use as an independent variable

Outputs:

- None

Notes:

- None

.....

```
self.data = data_df
self.y = data_df[[y_col]]
self.x = sm.add_constant(data_df[x_cols], prepend = False)
```

```
def fit_model(self, model_type, error_type = None):
    ....
```

Inputs:

- model_type: string indicating whether the model should be fit using OLS or a Robust Regression
- error_type: string indicating if model fit should use heteroscedasticity aware standard errors

- Note this should only be used if model_type = OLS; could add an error catch here to handle this

Outputs:

- model_summary: Dataframe containing the statsmodels regression output

```
=====
    """
if model_type == 'OLS':
    if error_type is not None:
        self.model = sm.OLS(self.y, self.x).fit(cov_type = error_type)
    else:
        self.model = sm.OLS(self.y, self.x).fit()

if model_type == 'Robust':
    self.model = sm.RLM(self.y, self.x).fit()

tmp_ = pd.DataFrame([self.model.fittedvalues, self.model.resid], index = ['yhat', 'e']).T
self.output_df = pd.merge(self.data, tmp_, left_index = True, right_index = True)

return self.model.summary2()

def plot_modelfit(self, model_name):
    """
Inputs:
    - model_name: string to title plot and identify model
    Outputs:
    - None; plot generated by call
    """
    sns.set_context('talk')
    sns.pairplot(data = self.output_df[['load', 'temp', 'e', 'hr']], hue = 'hr', height = 5, aspect = 2, palette = 'rocket')
    plt.suptitle(model_name)
    plt.tight_layout()
```

Data Processing

- Below I conduct some general data pre-processing, especially regarding peak vs. off-peak load
- We need (for latter questions especially) a way to designate on and off peak load
- I pulled NERC Holidays and added them to the assessment on the NERC_HOL sheet
 - Holiday dates for 2016-2019 were generated using
<https://www.energygps.com/HomeTools/PowerCalendar>
 - I checked total peak load hours using sums across month and year against the above link

In [6]:

```
df_raw = pd.read_excel('Data/Clean/TempLoadData.xlsx',
                      sheet_name='TempLoad',
                      usecols=['Hour Beginning', 'Temp (F)', 'Load (MW)'],
                      skiprows=1, parse_dates=True).rename(columns = {'Hour Beginning':'ts', 'Temp (F)':'temp', 'Load (MW)':'load'})

nerc_hol = pd.read_excel('Data/Clean/TempLoadData.xlsx',
                        sheet_name='NERC_HOL',
                        parse_dates=True).iloc[:-3]
```

In [7]:

```
df = gen_peak_id(df_raw, 'ts', nerc_hol, 'HolidayDate')
df.head()
```

Out[7]:

	ts	temp	load	yr	mth	dom	dow	hr	on_peak
0	2016-09-01 00:00:00	79.0	131.210245	2016	9	1	3	0	0.0
1	2016-09-01 01:00:00	79.0	123.057346	2016	9	1	3	1	0.0
2	2016-09-01 02:00:00	78.0	117.498482	2016	9	1	3	2	0.0
3	2016-09-01 03:00:00	77.0	113.287505	2016	9	1	3	3	0.0
4	2016-09-01 04:00:00	77.0	111.564578	2016	9	1	3	4	0.0

Exploratory Data Analysis

- I view EDA as equivalent to the axe sharpening in the alledged Lincoln quote "If I had eight hours to chop down a tree, I'd spend six sharpening my axe."
- Furthermore, I firmly believe if you cannot plot (or simulate) the process you are analyzing then you cannot properly model it; you don't understand it well enough

In [8]:

```
df_daily_7davg = df.set_index('ts').resample('D').mean().rolling(7).mean().dropna()
# Generate the daily average load, then create a rolling 7-day average to cut down on the noise for
# visualization's sake. Dropna() is called to get rid of first 6 days
df_daily_7davg.head()
```

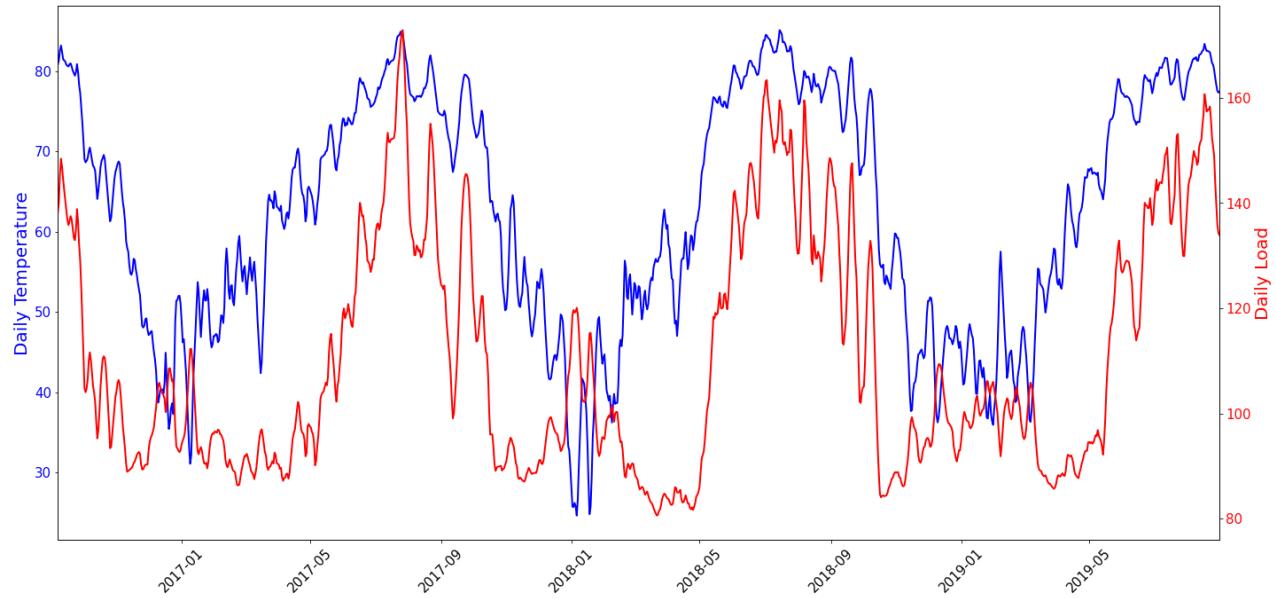
Out[8]:

	temp	load	yr	mth	dom	dow	hr	on_peak
ts								
2016-09-07	80.779762	138.181312	2016.0	9.0	4.0	3.0	11.5	0.380952
2016-09-08	81.345238	140.156468	2016.0	9.0	5.0	3.0	11.5	0.380952
2016-09-09	82.541667	145.302565	2016.0	9.0	6.0	3.0	11.5	0.380952
2016-09-10	83.160714	148.414973	2016.0	9.0	7.0	3.0	11.5	0.380952
2016-09-11	82.357143	146.854949	2016.0	9.0	8.0	3.0	11.5	0.380952

In [9]:

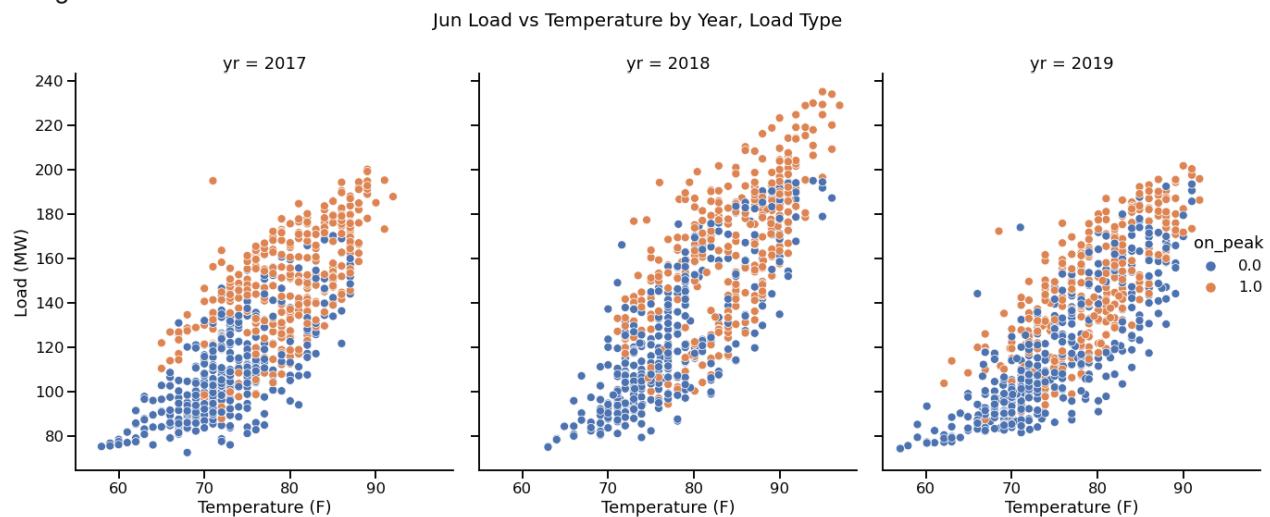
```
plot_dual_tseries(df_daily_7davg, ts_1='temp', ts_2='load',
                  ts_1_title='Daily Temperature', ts_2_title='Daily Load',
                  plot_title='Rolling 7 Day Daily Averages')
```

loadprediction_and_hedging
Rolling 7 Day Daily Averages



In [10]: `plot_annual_loadprofile(df, 6)`

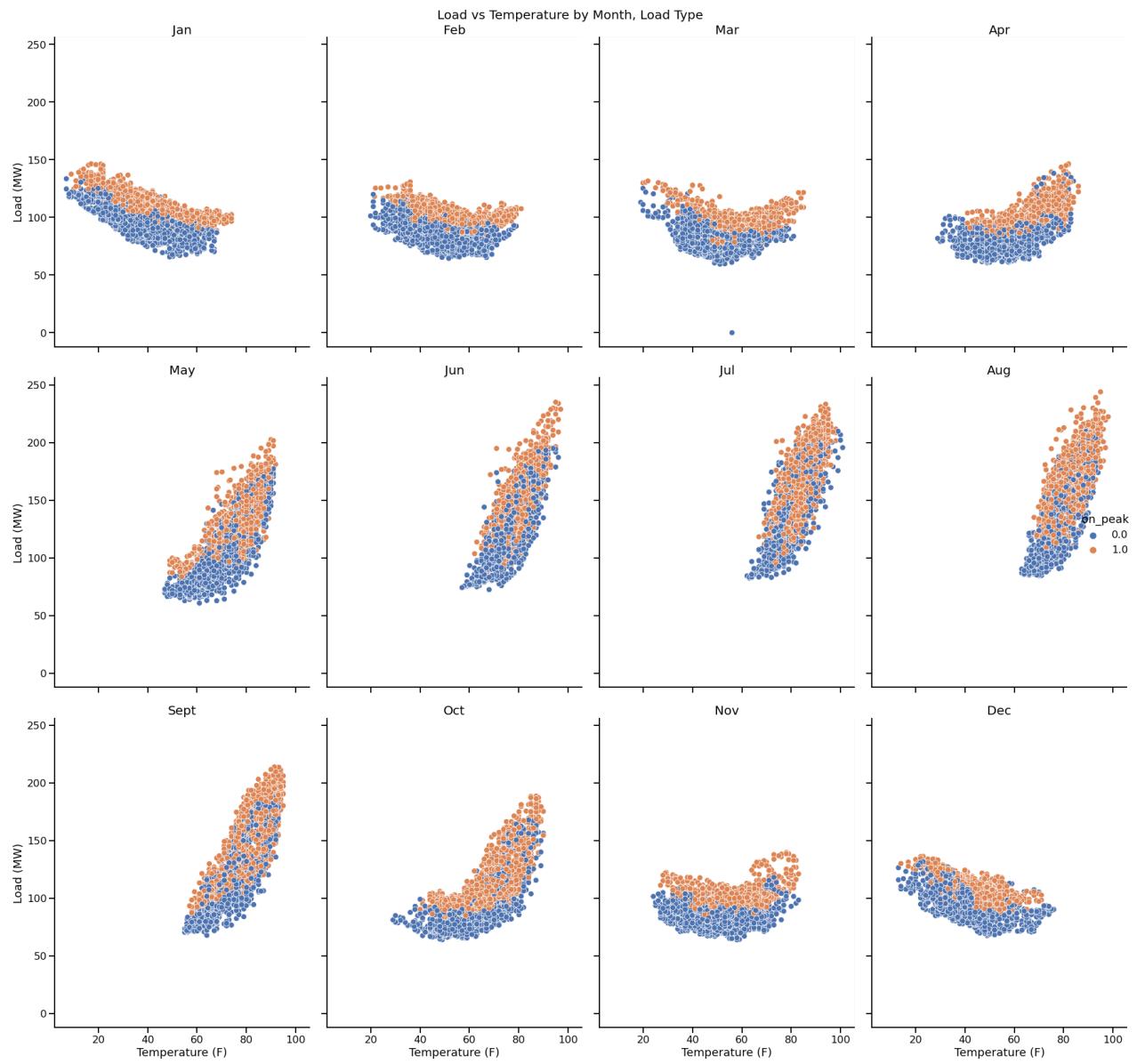
<Figure size 1440x720 with 0 Axes>



In [11]: `plot_monthly_loadprofile(df)`

<Figure size 1440x720 with 0 Axes>

loadprediction_and_hedging

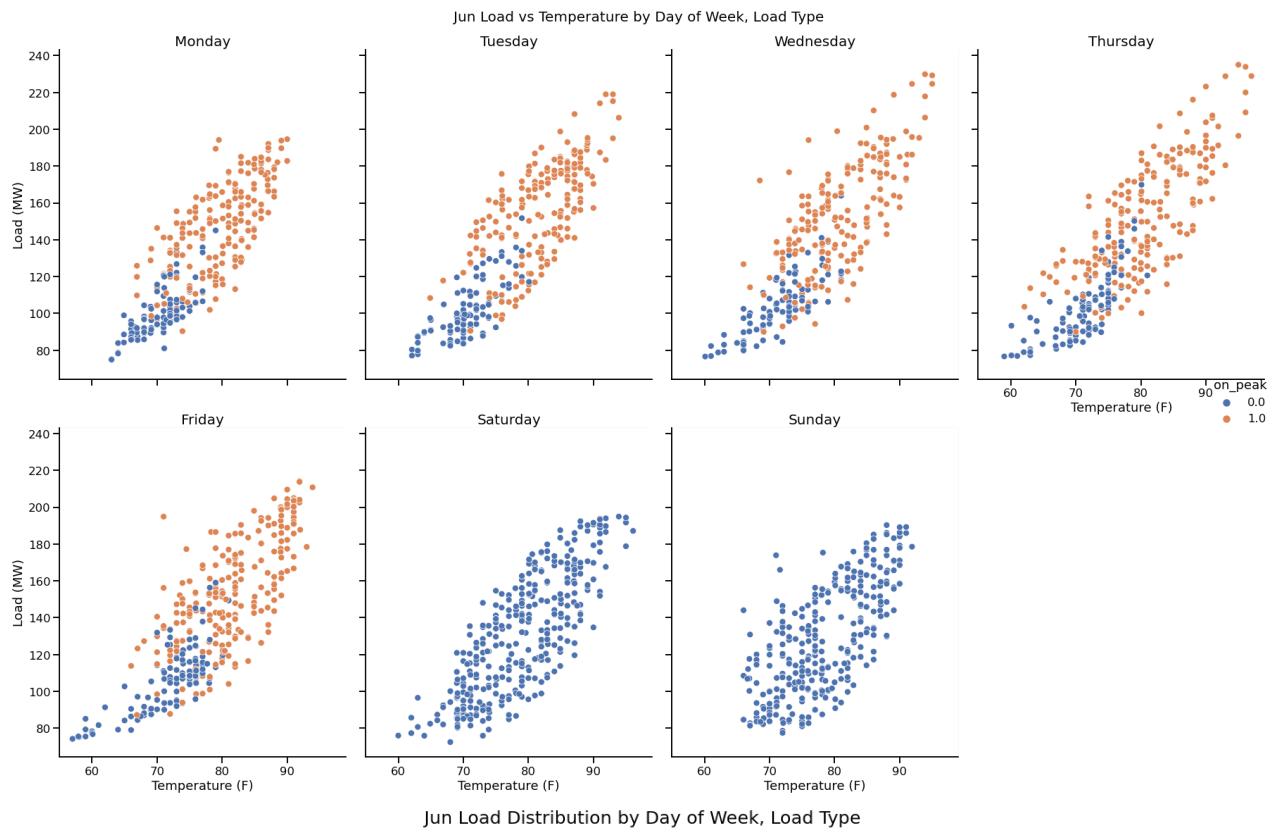


In [12]:

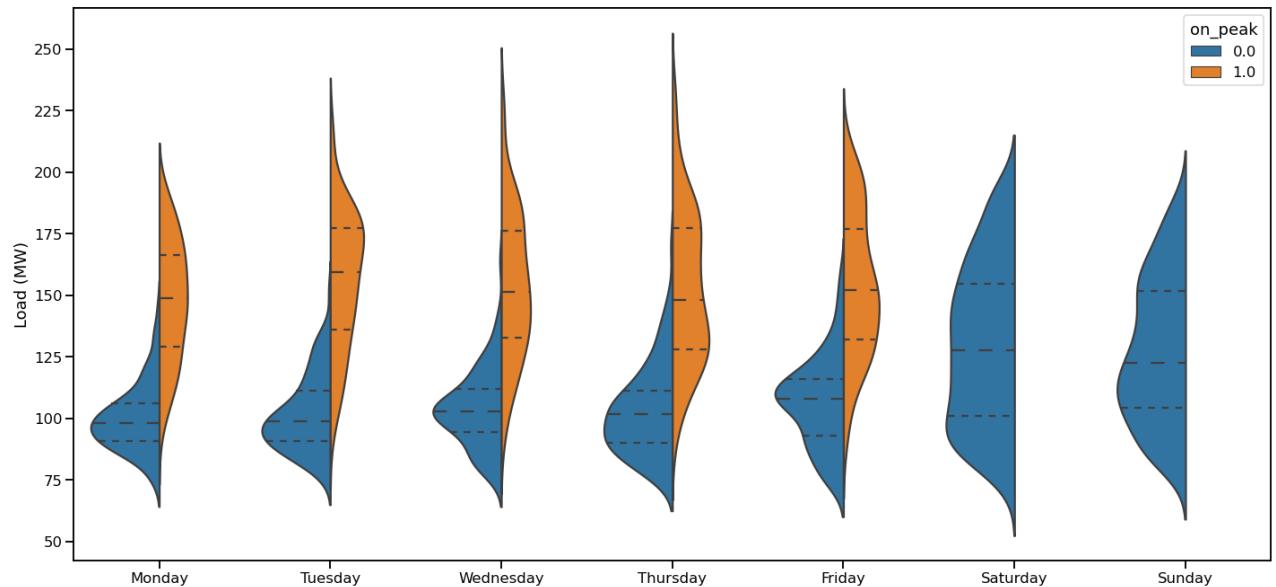
```
plot_dayofweek_loadprofile(df, 6)
```

<Figure size 1440x720 with 0 Axes>

loadprediction_and_hedging

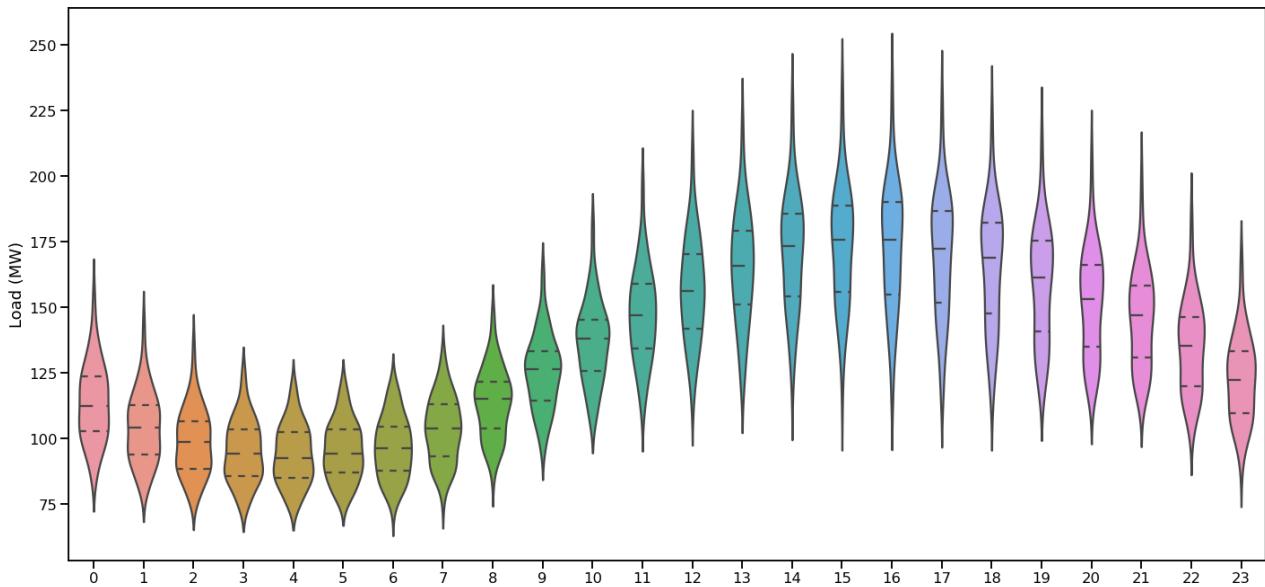


Jun Load Distribution by Day of Week, Load Type



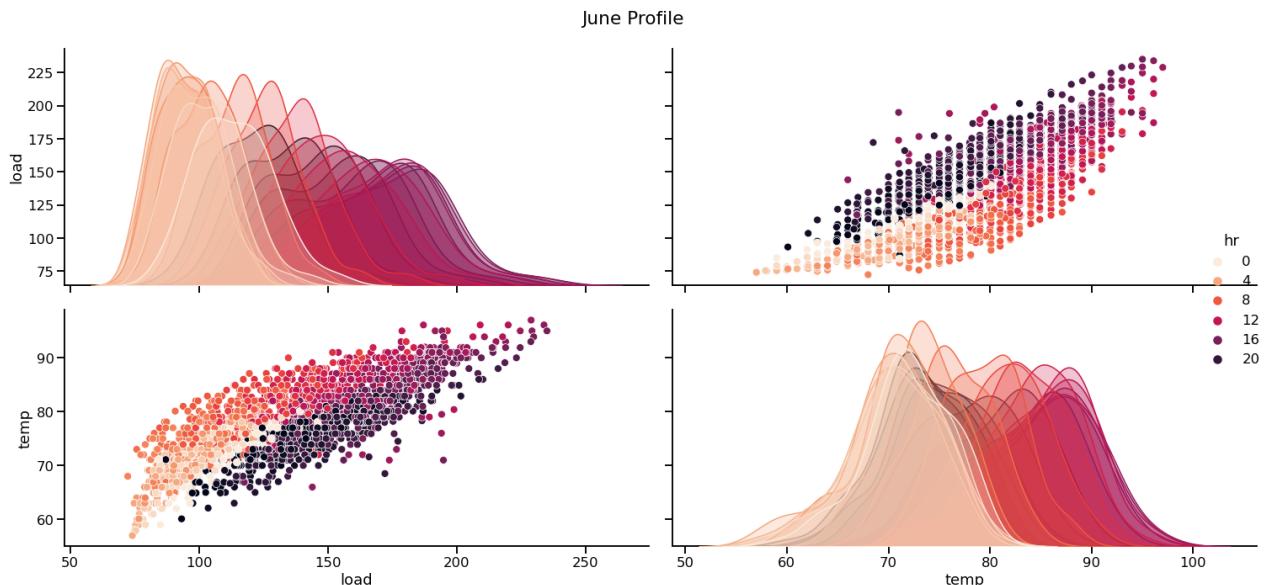
In [13]:

```
plot_hourly_loadprofile(df, 6)
```



In [14]:

```
df_j = df.loc[df['mth'] == 6] #only June obs.
sns.set_context('talk')
sns.pairplot(data = df_j[['load', 'temp', 'hr']], hue = 'hr', height = 5, aspect = 2, palette = 'rocket_r')
plt.suptitle('June Profile')
plt.tight_layout()
```



- Visualizing the load-temperature relationship annually, there does not appear to be significant variation while the opposite is true when viewing by month, day of week, and hour
 - Monthly load variation likely stems from natural seasonality while daily and hourly variation likely stem from the workweek and a workday/temperature cycle, respectively
- The hourly load profile displays a wave pattern which could be modeled with either a polynomial regression or sine wave for short term predictions
 - However, given we are forecasting at a longer-term horizon, I opted for a more structural model using dummy variables and interactions
 - Additionally, while fat tails (see next bullet) should give us some pause when using linear regression, polynomial regression is notoriously fragile when extrapolating out of sample

- Visualizing load by day of week and peak vs. off-peak, we can see both the mean and (more importantly) the distribution of load change markedly
 - Also note that our forecast day, a Thursday, displays the fattest tails
- Looking at the workweek, Monday's profile also differs markedly and thus I use Tuesday-Friday data to fit the model

Modeling

In [15]:

```
mask_1 = (df_j['dow'] > 0) & (df_j['dow'] < 5)
reg_df = df_j.loc[mask_1].drop(columns = ['yr', 'mth', 'dom', 'dow']).reset_index(drop = True)
reg_df.head()
```

Out[15]:

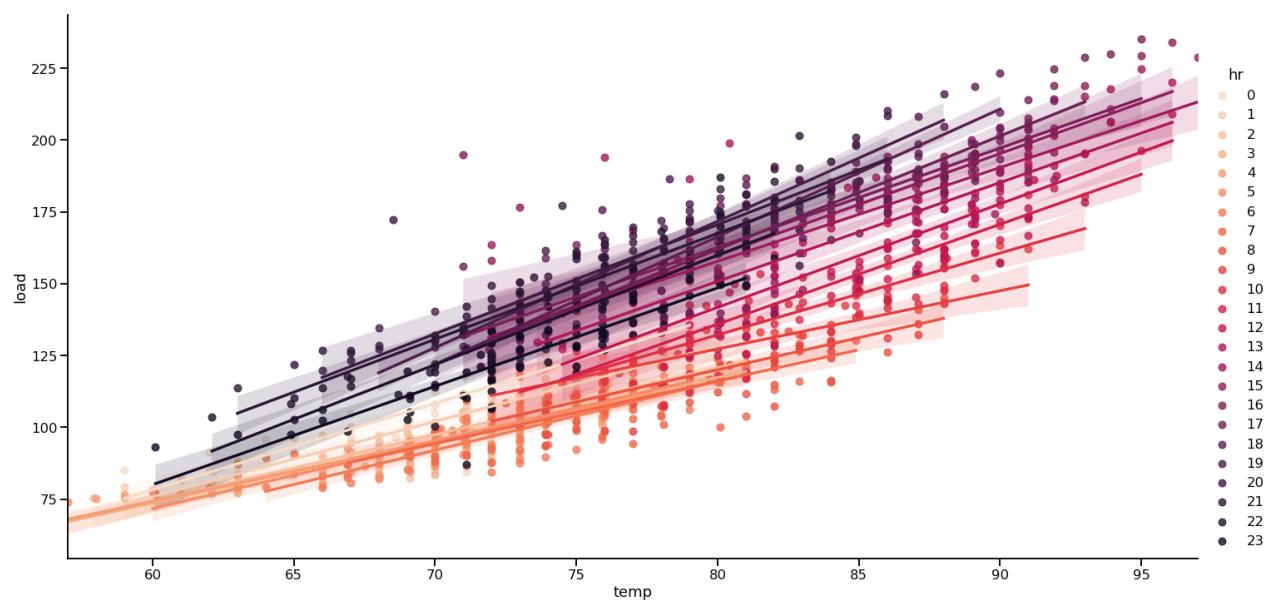
	ts	temp	load	hr	on_peak
0	2017-06-01 00:00:00	70.0	108.028450	0	0.0
1	2017-06-01 01:00:00	69.0	100.454888	1	0.0
2	2017-06-01 02:00:00	68.0	94.941490	2	0.0
3	2017-06-01 03:00:00	67.0	91.426984	3	0.0
4	2017-06-01 04:00:00	68.0	90.396514	4	0.0

- The below plot provides a hint to what we find as we iterate through models--the optimal model will likely require hour-specific intercepts AND slopes

In [16]:

```
sns.lmplot(y = 'load', x = 'temp', hue = 'hr', data = reg_df, height = 10, aspect = 2, palette = 'rocket_r')
```

Out[16]:



- We face two modeling choices

1. Whether to simply model differences in base load (intercepts) or also differences in load sensitivity (interactions) with the dummy variables
 2. What to use as dummy variables—peak vs. off-peak or hour dummies
- I opted for hourly dummies as there is an hourly pattern in the load distribution we would overlook if aggregated to just peak and off-peak dummies
 - All the models discussed below are fit to data from the Tues-Fri observations from June and utilize heteroscedasticity robust standard errors
 - Looking at the workweek using the daily visualizations in the prior section, Monday's profile also differs markedly and weekends are known to differ markedly as well

Model 1: $L_t = \alpha + \beta_1 F_T$

- I began with a basic model, $L=a+\beta_1 \text{Temp}$
 - This has the advantage of letting us look at the residuals by hour
 - There is clearly clustering in the residuals of the basic model which supports the dummy variable intercepts
 - The pattern between residuals and temperature (which would give insight regarding interactions) is less clear when looking at the basic model residuals
 - The Adjusted R2 is 0.695

```
In [17]: m1 = LoadModel(data_df = reg_df, y_col = 'load', x_cols = 'temp')
m1_results = m1.fit_model(model_type='OLS', error_type='HC1')
m1_results
```

C:\Users\Garrett\anaconda3\envs\tenaska_env\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWarning: In a future version of pandas all arguments of concat except for the argument 'objs' will be keyword-only
 x = pd.concat(x[:order], 1)

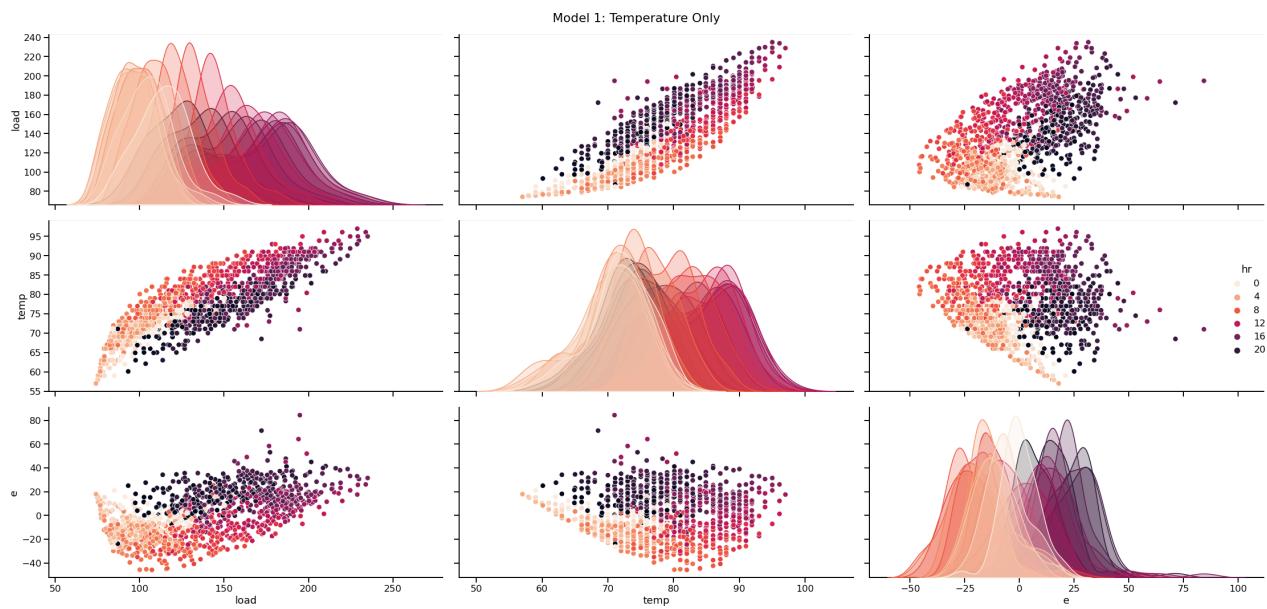
```
Out[17]:
```

	Model:	OLS	Adj. R-squared:	0.695		
Dependent Variable:	load		AIC:	10729.7419		
Date:	2021-07-19 03:20		BIC:	10739.9617		
No. Observations:	1224		Log-Likelihood:	-5362.9		
Df Model:	1		F-statistic:	3644.		
Df Residuals:	1222		Prob (F-statistic):	0.00		
R-squared:	0.695		Scale:	374.88		
	Coef.	Std.Err.	z	P> z 	[0.025	0.975]
temp	3.8760	0.0642	60.3631	0.0000	3.7501	4.0019
const	-164.6876	4.9825	-33.0534	0.0000	-174.4531	-154.9222

Omnibus:	19.673	Durbin-Watson:	0.195
Prob(Omnibus):	0.000	Jarque-Bera (JB):	18.354
Skew:	0.257	Prob(JB):	0.000
Kurtosis:	2.689	Condition No.:	813

In [18]:

```
m1.plot_modelfit('Model 1: Temperature Only')
```



$$\text{Model 2: } L_t = \alpha + \sum_{i=0}^{23} \{\lambda_i H_i\} + \beta_1 F_T$$

In [19]:

```
reg2_df = reg_df.copy()
hr_dummies = pd.get_dummies(reg2_df['hr'], drop_first=True, prefix='H')
reg2_df = reg2_df.merge(hr_dummies, left_index=True, right_index=True)
reg2_df.head()
```

Out[19]:

	ts	temp	load	hr	on_peak	H_1	H_2	H_3	H_4	H_5	...	H_14	H_15	H_16	H_17	I
0	2017-06-01 00:00:00	70.0	108.028450	0		0.0	0	0	0	0	...	0	0	0	0	0
1	2017-06-01 01:00:00	69.0	100.454888	1		0.0	1	0	0	0	...	0	0	0	0	0
2	2017-06-01 02:00:00	68.0	94.941490	2		0.0	0	1	0	0	...	0	0	0	0	0
3	2017-06-01 03:00:00	67.0	91.426984	3		0.0	0	0	1	0	...	0	0	0	0	0
4	2017-06-01 04:00:00	68.0	90.396514	4		0.0	0	0	0	1	0	...	0	0	0	0

5 rows × 28 columns

In [20]:

loadprediction_and_hedging

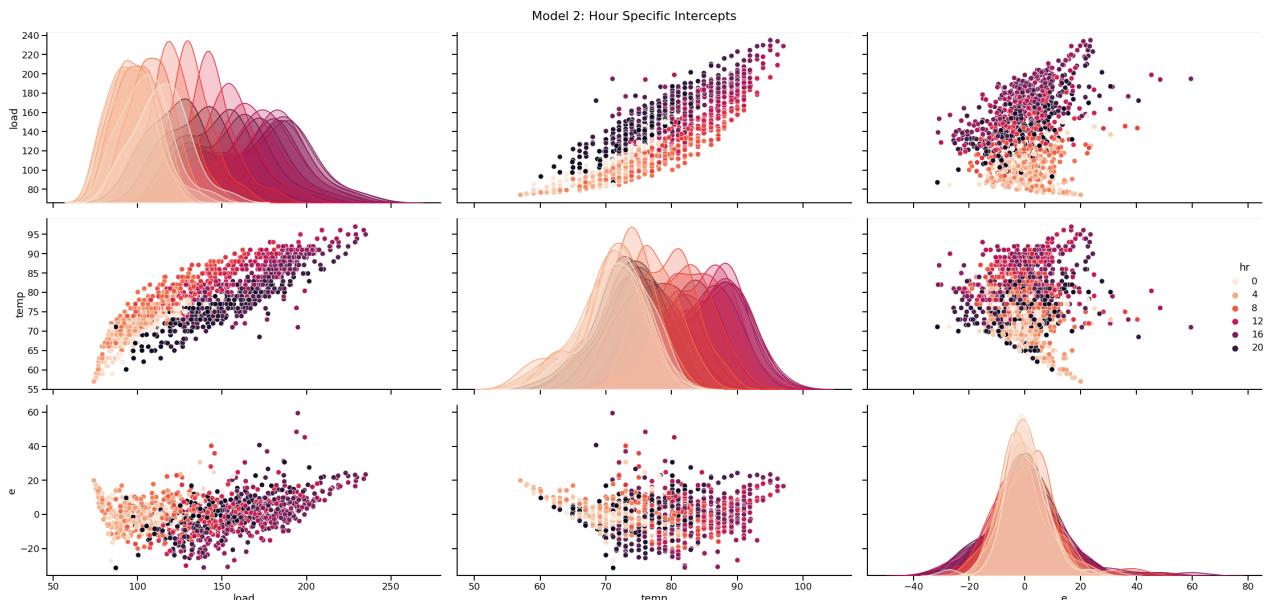
H_23	6.9243	1.6555	4.1825	0.0000	3.6795	10.1691
temp	3.1729	0.0812	39.0833	0.0000	3.0137	3.3320
const	-114.0204	5.9171	-19.2698	0.0000	-125.6177	-102.4232

Omnibus: 108.252 Durbin-Watson: 0.366
 Prob(Omnibus): 0.000 Jarque-Bera (JB): 302.268
 Skew: 0.460 Prob(JB): 0.000
 Kurtosis: 5.254 Condition No.: 1945

- We could conduct a Lagrange Multiplier test to test Model 2 vs. Model 1, but given all the dummies are significant we know Model 2 offers a superior fit

In [21]:

```
m2.plot_modelfit('Model 2: Hour Specific Intercepts')
```



$$\text{Model 3: } L = \alpha + \sum_{i=0}^{23} (\lambda_i H_i) + \beta_1 F + \sum_{i=0}^{23} (\phi_i H_i F_t)$$

- To test whether the intercepts and interactions are significant, I conducted a joint Likelihood Multiplier test of the intercept-interaction model vs. the intercept model
 - This stepwise process (thinking of our iteration from Model 1 to 2 to 3) is valid because the models are nested
- The LM test returns a p-value < 0.01, which supports the least restricted model
 - By least restricted I mean the model which includes both hourly intercepts and hourly interactions with temperature.
 - Thus, each hour has a distinct base load and sensitivity to temperature: $L = a + \sum \lambda_i H_i + \beta_1 Temp + \sum \phi_i H_i Temp$

Skew: 0.606 Prob(JB): 0.000
 Kurtosis: 6.634 Condition No.: 30560

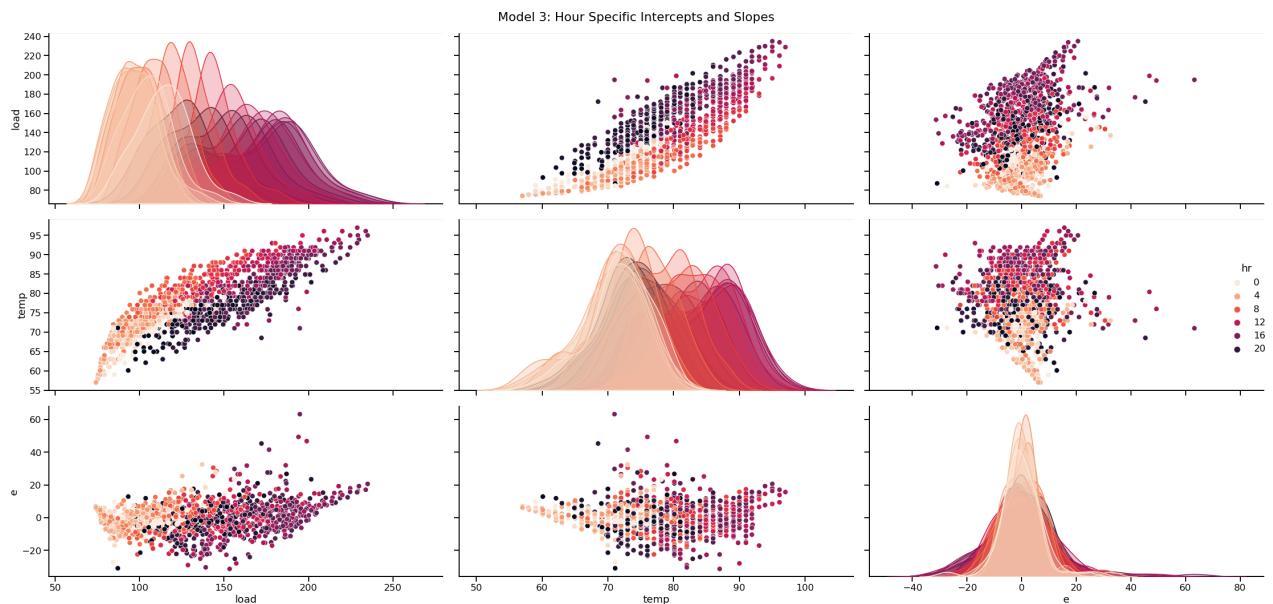
- Given the significance of the slope terms is mixed, I conduct a LM test of overall model significance vs. Model 2; the results below support Model 3

In [24]: `m3.model.compare_lm_test(m2.model)`

```
C:\Users\Garrett\anaconda3\envs\tenaska_env\lib\site-packages\statsmodels\regression\linear_model.py:20
86: FutureWarning: Support for multi-dimensional indexing (e.g. `obj[:, None]`) is deprecated and will be removed in a future version. Convert to a numpy array before indexing instead.
```

scores = wexog * wresid[:, None]
 Out[24]: (142.29949120575606, 3.4713934404480327e-19, 23.0)

In [25]: `m3.plot_modelfit('Model 3: Hour Specific Intercepts and Slopes')`



- Given the presence of outliers, especially on our prediction day of Thursday, I chose to implement a robust regression version of Model 3 for the final prediction
 - This minimizes the mean absolute deviation instead of squared deviations
 - MAD is more robust in the presence of fat tails
- I arrived at the predictions by averaging model fitted values across all Thursdays in June

In [26]: `robust_m3 = LoadModel(data_df = reg3_df, y_col = 'load', x_cols = list(reg3_df.columns[5:])+['temp'])
 robust_m3_results = robust_m3.fit_model(model_type='Robust')
 robust_m3_results`

```
C:\Users\Garrett\anaconda3\envs\tenaska_env\lib\site-packages\statsmodels\tsa\tsatools.py:142: FutureWar
```


Discussion:

- Robust regression will help mitigate the impact but will not cure it. However, the residuals look relatively Gaussian.
- Given weather is typically modeled as a normal distribution, I would seek to improve the model by looking for alternative drivers behind the fat tails.
- The forecast would also improve as the forecast horizon approaches.

In [30]:

```

df_jth = df_j.loc[df_j['dow'] == 3]
df_jth['ts'] = pd.DatetimeIndex(df_jth['ts'], normalize=True)
df_jth.set_index('ts', inplace = True)

june_thurs = ['2017-06-01', '2017-06-08', '2017-06-15', '2017-06-22', '2017-06-29',
              '2018-06-07', '2018-06-14', '2018-06-21', '2018-06-28',
              '2019-06-06', '2019-06-13', '2019-06-20', '2019-06-27']

df_jth_maxdf = pd.DataFrame([[df_jth.loc[dt, 'load'].idxmax().hour, df_jth.loc[dt, 'load'].max(),
                               df_jth.loc[dt, 'temp'].idxmax().hour, df_jth.loc[dt, 'temp'].max()] for dt in june_thurs],
                             columns=['maxload_hr', 'maxload', 'maxtemp_hr', 'maxtemp'], index = june_thurs)

sns.set_context('talk')
g2 = sns.jointplot(x = 'maxload_hr', y = 'maxload', data = df_jth_maxdf,
                    kind = 'kde', fill = True, height = 10, palette = 'rocket_r')
g2.fig.suptitle('Load: Peak Value and Peak Hour', fontsize = 20)
g2.set_axis_labels('Peak Hour', 'Max Load')
g2.fig.tight_layout()

```

Load: Peak Value and Peak Hour

