# Analyze customer data from Austin, TX.

In this lab, we will analyze various features of customer energy usage data from the Pecan Street dataset. The dataset we aim to explore is for the month of January 2017 from various customers in Austin, TX.

Distinguish which houses have rooftop photovoltaic (PV) panels and which ones do not. For the houses that do not have PV panels, try to predict how much money they would save, given Austin's retail energy rates.

Given the temperature profile, can you predict which houses have AC?

```python
# Start with customary imports.

import tensorflow as tf
import numpy as np
import pandas as pd
import random
import datetime
import math
from sklearn.model_selection import train_test_split
from sklearn import preprocessing, linear_model
from sklearn.metrics import classification_report
import matplotlib.pyplot as plt
```

## Understand and parse the dataset

There are two csv files that we will utilize in this lab:

1. 'dataport-metadata.csv': It contains the details of what data is available from each house in the entire Pecan Street dataset.
2. 'July2017.csv': It contains the energy data from houses in Austin for the month of July.

### Parse the metadata file.

Only consider customers from Austin, TX for which electricity usage data is present.

```python
In [ ]: # Load the 'metadata' file that contains information about individual houses i
        dfCityCustomers = pd.read_csv("dataport-metadata.csv", index_col=0)

        # Only consider the houses that are in Austin and have power consumption data.
        dfCityCustomers = dfCityCustomers.loc[(dfCityCustomers['city'] == 'Austin') &

        # Restrict attention to useful columns.
        dfCityCustomers = dfCityCustomers[['date_enrolled', 'date_withdrawn',
                                            'building_type', 'total_square_footage', 'f
                                            'pv', 'air1', 'air2', 'air3', 'airwindowuni
                                            'gen', 'use', 'grid']]

        # Replace binary data with zeros and ones.
        binaryColumns = ['pv', 'air1', 'air2', 'air3', 'airwindowunit1', 'gen', 'use',
        for bColumn in binaryColumns:
            dfCityCustomers[bColumn] = dfCityCustomers[bColumn].map({'yes' : 1}).filln

        start_day = datetime.datetime.strptime('2017-07-01', '%Y-%m-%d').date()

        dfCityCustomers['date_enrolled'] = [datetime.datetime.strptime(x, "%Y-%m-%d")
                                            for x in dfCityCustomers['date_enrolled']]

        dfCityCustomers = dfCityCustomers.loc[(dfCityCustomers['date_enrolled'] <= '20
        print("Parsed the metadata file successfully.")
```

Parsed the metadata file successfully.


## Parse the energy usage data.

Here, we shall clean the data.

1. The data downloaded from dataport.cloud has a mislabeled column. Correct that.
2. Split the 'localhour' field into two fiels: actual date and an hour of day.
3. Only choose data from households that are 'Single-Family Homes'.
4. Make sure there is data from 31 days.
5. Ensure that the metadata includes the square footage for the entire home and its first floor.

Finally, create a pandas dataframe where indices are house id's.

In [5]:
```python
# Load the data from Jan 2017 from houses in Texas.
dfData = pd.read_csv("July2017.csv")

# Cleanup item 1. Alter the column names because there is an error in the down
# 'dataid' and 'localhour' are switched in the .csv file. Correct it.
# Cleanup item 2. Split the hour and the day in the field 'localhour' and conv
# actual datetime object.

dfData.columns = [x.lstrip() for x in dfData.columns]
dfData = dfData[['dataid', 'localhour', 'use', 'grid', 'gen', 'air1', 'air2',
dfData.columns = ['date', 'dataid', 'use', 'grid', 'gen', 'air1', 'air2', 'air

dfData['hour'] = [datetime.datetime.strptime(x[:-3], "%Y-%m-%d %H:%M:%S").hour
dfData['date'] = [datetime.datetime.strptime(x[:-3], "%Y-%m-%d %H:%M:%S").date

# Create a dataframe where indices are house id's.
dfData_houses = pd.DataFrame(columns=['ac', 'pv', 'area', 'area_floor', 'data'

for house_id in dfData['dataid'].unique():

    # Make sure that each house with consumption data is also in metadata.
    if house_id in dfCityCustomers.index.values:
        dfData_one_house = dfData.loc[dfData['dataid'] == house_id]

        # Cleanup item 3, 4, 5.
        if ((len(dfData_one_house) >= 24 * 31) &
        (np.sum(dfData_one_house['use'].values) != 0) &
        (dfCityCustomers.loc[house_id, 'building_type'] == 'Single-Family Home
        (np.isnan(dfCityCustomers.loc[house_id, 'total_square_footage']) == Fa
        (np.isnan(dfCityCustomers.loc[house_id, 'first_floor_square_footage'])

            # Create a pandas dataframe with house id's as indices and has the
            # 1. Binary status: "ac", "pv".
            # 2. Total square footrage, and the square footage of the first fl

            is_ac_in_house =  (dfCityCustomers.loc[house_id, 'air1']
                               or dfCityCustomers.loc[house_id, 'air2']
                               or dfCityCustomers.loc[house_id, 'air3']
                               or dfCityCustomers.loc[house_id, 'airwindowunit
                              )

            dfData_houses.loc[house_id] = [is_ac_in_house,
                                           dfCityCustomers.loc[house_id, 'pv']
                                           dfCityCustomers.loc[house_id, 'tota
                                           dfCityCustomers.loc[house_id, 'firs
                                           dfData_one_house
                                          ]

# Define a function that retrieves the hourly energy profile from the column i
# specific house and day. The day is measured as number of days since the star

def energy_day(house_id, day, field):
    date_day = start_day + datetime.timedelta(days=day)
    dfData_house = dfData_houses.loc[house_id, 'data']
    return dfData_house.loc[dfData_house['date'] == date_day].sort_values(by=[

# Define a function that retrieves the aggregate energy consumed (or produced)
```

```python
# named "field" from a specific house over all days.

def total_energy_all_days(house_id, field):
    total_energy = 0
    for day in range(31):
        total_energy += np.sum(energy_day(house_id, day, field))
    return total_energy


# Cleanup item 4 continued: Delete data from houses where the date and hour do
# align with the 31 days in July, and hours being from 0 to 23.
house_ids_to_delete = []
for house_id in dfData_houses.index.values:
    for day in range(31):
        if len(energy_day(house_id, day, 'grid')) != 24:
            house_ids_to_delete.append(house_id)
            continue

dfData_houses.drop(house_ids_to_delete, inplace=True)

print("Data loaded and parsed successfully from %d single-family homes." % (le

del dfData , dfCityCustomers
```

```
---------------------------------------------------------------------------
ParserError                               Traceback (most recent call last)
<ipython-input-5-0e431a09ac40> in <cell line: 2>()
      1 # Load the data from Jan 2017 from houses in Texas.
----> 2 dfData = pd.read_csv("July2017.csv")
      3
      4 # Cleanup item 1. Alter the column names because there is an error in
the downloaded data. The column names
      5 # 'dataid' and 'localhour' are switched in the .csv file. Correct it.

/usr/local/lib/python3.10/dist-packages/pandas/util/_decorators.py in wrapper
(*args, **kwargs)
    209                 else:
    210                     kwargs[new_arg_name] = new_arg_value
--> 211             return func(*args, **kwargs)
    212
    213         return cast(F, wrapper)

/usr/local/lib/python3.10/dist-packages/pandas/util/_decorators.py in wrapper
(*args, **kwargs)
    329                 stacklevel=find_stack_level(),
    330             )
--> 331             return func(*args, **kwargs)
    332
    333         # error: "Callable[[VarArg(Any), KwArg(Any)], Any]" has no

/usr/local/lib/python3.10/dist-packages/pandas/io/parsers/readers.py in read_
csv(filepath_or_buffer, sep, delimiter, header, names, index_col, usecols, sq
ueeze, prefix, mangle_dupe_cols, dtype, engine, converters, true_values, fals
e_values, skipinitialspace, skiprows, skipfooter, nrows, na_values, keep_defa
ult_na, na_filter, verbose, skip_blank_lines, parse_dates, infer_datetime_for
mat, keep_date_col, date_parser, dayfirst, cache_dates, iterator, chunksize,
compression, thousands, decimal, lineterminator, quotechar, quoting, doublequ
ote, escapechar, comment, encoding, encoding_errors, dialect, error_bad_line
s, warn_bad_lines, on_bad_lines, delim_whitespace, low_memory, memory_map, fl
oat_precision, storage_options)
    948     kwds.update(kwds_defaults)
    949
--> 950     return _read(filepath_or_buffer, kwds)
    951
    952

/usr/local/lib/python3.10/dist-packages/pandas/io/parsers/readers.py in _read
(filepath_or_buffer, kwds)
    609
    610     with parser:
--> 611         return parser.read(nrows)
    612
    613

/usr/local/lib/python3.10/dist-packages/pandas/io/parsers/readers.py in read
(self, nrows)
   1776                 columns,
   1777                 col_dict,
-> 1778             ) = self._engine.read(  # type: ignore[attr-defined]
   1779                 nrows
   1780             )
```

```
/usr/local/lib/python3.10/dist-packages/pandas/io/parsers/c_parser_wrapper.py
in read(self, nrows)
    228            try:
    229                if self.low_memory:
--> 230                    chunks = self._reader.read_low_memory(nrows)
    231                    # destructive to chunks
    232                    data = _concatenate_chunks(chunks)

/usr/local/lib/python3.10/dist-packages/pandas/_libs/parsers.pyx in pandas._l
ibs.parsers.TextReader.read_low_memory()

/usr/local/lib/python3.10/dist-packages/pandas/_libs/parsers.pyx in pandas._l
ibs.parsers.TextReader._read_rows()

/usr/local/lib/python3.10/dist-packages/pandas/_libs/parsers.pyx in pandas._l
ibs.parsers.TextReader._tokenize_rows()

/usr/local/lib/python3.10/dist-packages/pandas/_libs/parsers.pyx in pandas._l
ibs.parsers.raise_parser_error()

ParserError: Error tokenizing data. C error: Expected 71 fields in line 2570,
saw 77
```

# Distinguish houses with rooftop solar panels from daily energy usage profile.

Take data of energy drawn from the grid for 10 days and do logistic regression.

## Q1. Explain Logistic regression. (10 points)

Logic Regression is a classification method that determines its output based on one or more inputs. Its outcome is discrete (0 or 1) and it is put through a logistic function that combines the inputs.

```
In [ ]:  print("Number of houses with PV panels = %d" % (len(dfData_houses.loc[dfData_h
         print("Number of houses without PV panels = %d" % (len(dfData_houses.loc[dfDat

         XX = []
         YY = []

         days_data = random.sample(range(31), 10)

         for house_id in dfData_houses.index.values:
             XX.append(np.ravel([energy_day(house_id, day, 'grid') for day in days_data
             YY.append(dfData_houses.loc[house_id, 'pv'])

         YY = np.reshape(YY, (-1, 1))

         train_X, test_X, train_Y, test_Y = train_test_split(XX, YY, test_size=.2, shuf

         train_X = tf.dtypes.cast(train_X, tf.float32)
         test_X = tf.dtypes.cast(test_X, tf.float32)
         train_Y = tf.dtypes.cast(train_Y, tf.float32)
         test_Y = tf.dtypes.cast(test_Y, tf.float32)

         del XX, YY
```

```
In [ ]:  print("Number of houses with PV panels = %d" % (len(dfData_houses.loc[dfData_h
         print("Number of houses without PV panels = %d" % (len(dfData_houses.loc[dfDat
```

## Q2. Design the neural network. In the next cell, fill in the missing pieces. (30 points)

```python
In [ ]: nDimX = np.shape(train_X)[1]
        nDimY = np.shape(train_Y)[1]

        weight = # Enter code here
        bias = # Enter code here

        trainable_variables = [weight, bias]

        @tf.function
        def neuralNetworkModel(X):
            global weight, bias
            # enter code here (hint: tf.nn.sigmoid may be useful)

        loss = # Enter code here (hint: tf.losses.BinaryCrossentropy may be useful)
        optimizer = # Enter code here (hint: the learning rate may need to be small)

        # Define number of epochs
        nEpoch = 1000

        # Define the training scheme
        def train(model, x_set, y_set):
            for epoch in range(nEpoch):
                # Fit the data and compute the gradients
                with tf.GradientTape() as tape:
                    prediction = model(x_set)
                    loss = loss_fn(y_true=y_set, y_pred=prediction)

                    # Print update
                    lossEpoch = loss.numpy()
                    print("Epoch: %d, Loss: = %1.1f" % (epoch + 1, lossEpoch))

                    # Optimize the weights
                    gradients = tape.gradient(loss, trainable_variables)
                    optimizer.apply_gradients(zip(gradients, trainable_variables))


        # Train the model
        print ("Start neural network training.")
        train(neuralNetworkModel, train_X, train_Y)

        test_prediction = tf.math.round(neuralNetworkModel(test_X))
        test_accuracy = tf.math.reduce_mean(tf.dtypes.cast(tf.math.equal(test_Y, test_
        print("Accuracy of logistic regression on test data = %.2f percent." % (test_a
```

### Q3. Print the classification report on the test data. The function 'classification_report' from 'sklearn.metrics' might prove useful. (10 points)

In [ ]:
```python
# Enter code here
```

### Q4. Based on the classification report you obtain, your classifier is better in which of the following tasks? (20 points)

1. If it identifies a house to have a PV panel, then it has a PV panel.
2. If there is a PV panel, then it identifies that it has a PV panel.

Furthermore, complete the code below to plot the energy drawn from the grid from houses with and without PV panels.

The classifier is better at identifying the second observation due to its simplicity.

In [ ]:
```python
# Plot energy drawn from grid for houses with PV's.
house_id_pv = random.sample(list(dfData_houses.loc[dfData_houses['pv'] == 1].i

fig, axs = plt.subplots(1, 5, sharey=True, figsize=(15,5))

for tt, house_id in enumerate(house_id_pv):
    for day in days_data:
        # Enter code here
        axs[tt].set_title("House " + str(house_id))
fig.suptitle('Houses with PV panels.', fontsize=18)

for ax in axs.flat:
    ax.set(xlabel='Hour', ylabel='Energy in kWh')
for ax in axs.flat:
    ax.label_outer()

house_id_not_pv = random.sample(list(dfData_houses.loc[dfData_houses['pv'] ==

fig, axs = plt.subplots(1, 5, sharey=True, figsize=(15,5))
axs = axs.ravel()

for tt, house_id in enumerate(house_id_not_pv):
    for day in days_data:
        # Enter code here
        axs[tt].set_title(house_id)

fig.suptitle('Houses without PV panels.', fontsize=18)

for ax in axs.flat:
    ax.set(xlabel='Hour', ylabel='Energy in kWh')
for ax in axs.flat:
    ax.label_outer()
```

**Q5. Design a classifier to distinguish between houses with and without PV based on the plots. Print your classification report, and compare it with logistic regression. (10 points, bonus)**

In [ ]:
```
# Compare the performance of a neural network based classifier with an educate
```

## What appliance consumes the most power?

Thermal loads are almost always the appliances that consume the most power. The amount of power they draw also typically grows with the size of the house. In the following, we have two tasks:

1. What percentage of energy consumption is due to an air conditioner?
2. Can you derive a linear relationship between household square footage and power consumption from air conditions in July 2017?

In [ ]:
```
for house_id in dfData_houses.index.values:
    dfData_houses.at[house_id, 'ac_usage'] = (total_energy_all_days(house_id,
                                               total_energy_all_days(house_id,
                                               total_energy_all_days(house_id,
                                               total_energy_all_days(house_id,
                                              )
    dfData_houses.at[house_id, 'total_usage'] = total_energy_all_days(house_id

fig, axs = plt.subplots(1, 2, figsize=(15,5))

axs[0].scatter(dfData_houses['area'].values,
               dfData_houses['ac_usage'].values,
               c='r', marker='o', label='AC usage'
              )
axs[0].set_xlabel('Floor area of house (sq. ft.)')
axs[0].set_ylabel('Total AC Usage (kWh).')
axs[0].set_title('Power consumption from air conditioners July 2017.', fontsiz

axs[1].hist(x=np.divide(dfData_houses['ac_usage'], dfData_houses['total_usage'
axs[1].set_xlabel('Fraction of AC usage over total power consumption')
axs[1].set_ylabel('Frequency')
axs[1].set_title('Histogram of AC usage as a fraction of total power consumpti
```

**Q6. Based on the above analysis, will you expect the total power consumption from customers to be more or less in October as compared to that in July? How does the above analysis compare to your analysis in Lab 1 on aggregate load prediction in Texas? (10 points)**

**Q7. Can you think of a business case for the above histogram? (5 points, bonus)**

6. I expect the power consumption to be less in October than July. It compares similarly to Lab 1 due to how the lower power needed to cool the houses after the summer is over.
7. A business case that could be used for the data is determining power consumption needed throughout a summer. This data can be good for estimating cost and consumption of power that businesses may require.

## Compute monthly electric bills for each customer

Your electricity bill consists of various charges. These charges depend on the utility company you pay your electricity bill to. In Champaign, IL, your distribution utility company is Ameren. In Austin, a major distribution utility company is Austin Energy. Their bill structure is discussed in the following link: https://austinenergy.com/ae/residential/rates/residential-electric-rates-and-line-items (https://austinenergy.com/ae/residential/rates/residential-electric-rates-and-line-items)

Calculate the monthly bill of each household.

1. For each customer, compute the total energy consumed, available in the data field "use" over the month. Use the tiered rate structure to compute the total power bill for energy consumption.
2. For customers with PV panels, compute the total energy produced by the PV panels. Assume that Ameren Energy pays 9.7 cents/kWh for such production, and subtract the amount for solar power production from the power bill.

```
In [ ]:  # Define a function that computes the electricity bill according to the struct

         def electricity_bill(consumption):
             customer_charge = 10
             power_supply_adjustment = 2.895 * consumption / 100.0
             community_benefit_charge = (0.154 + 0.124 + 0.335) * consumption / 100.0
             regulatory_charge = 1.342 * consumption / 100.0

             tier_rate = [2.801, 5.832, 7.814, 9.314, 10.814]
             tier_limits = [0, 500, 1000, 1500, 2500, math.inf]
             n_tiers = 5

             energy_charge = 0

             for tier in range(n_tiers):
                 consumption_tier = min(max(consumption, tier_limits[tier]), tier_limit
                 energy_charge += consumption_tier * tier_rate[tier] / 100.0

             return float('%.2f'%(1.01 * (energy_charge + customer_charge + power_suppl
                         + regulatory_charge)))

         for house_id in dfData_houses.index.values:
             dfData_houses.at[house_id, 'consumption_bill'] = electricity_bill(dfData_h
             dfData_houses.at[house_id, 'pv_savings'] = total_energy_all_days(house_id,
             dfData_houses.at[house_id, 'electricity_bill'] = dfData_houses.loc[house_i
             - dfData_houses.loc[house_id, 'pv_savings']


         print("Electricity bill computed for all customers.")
```

## What is a good indicator of electricity bill for consumption and savings from PV?

Electricity consumption significantly depends on floor area of a house. The dependency is even stronger, if the house is equipped with central AC. A scatter plot of the consumption bill against floor area reveals this dependency.

Monetary savings from PV panels depends on how many PV panels there are, which way they face, and how that relates to solar insolation. The number of panels installed largely depends on the roof area. The floor area of one of the floors is a good indicator.

```python
In [ ]: fig, axs = plt.subplots(1, 2, figsize=(15,4))

        # Plot of electricity bill for consumption against the square footage of the h

        axs[0].scatter(dfData_houses['area'].values,
                       dfData_houses['consumption_bill'].values,
                       c='r', marker='o', label='Electricity bill from consumption'
                  )
        axs[0].set_xlabel('Floor area of house (sq. ft.)')
        axs[0].set_ylabel('Energy bill (dollars).')
        axs[0].set_title('Electricity bill from power consumption over July 2017.', fo

        # Plot monetary savings from PV against the area of the first floor.

        houses_without_pv = dfData_houses.loc[dfData_houses['pv'] == 0].index.values
        houses_with_pv = dfData_houses.loc[dfData_houses['pv'] == 1].index.values

        axs[1].scatter(dfData_houses.loc[houses_with_pv, 'area'].values,
                       dfData_houses.loc[houses_with_pv, 'pv_savings'].values,
                       c='g', marker='o', label='Savings from PV')

        axs[1].set_xlabel('Area of first floor (sq. ft.)')
        axs[1].set_ylabel('Monthly savings (dollars).')
        axs[1].set_title('Monthly savings from PV panels over July 2017.', fontsize=14
```

## Given the square footage of the house and the first floor, compute the electricity bill with and without PV.

Perform a linear regression on electricity bill and monthly savings from PV. Use the linear fits to compute the anticipated electricity bill with and without PV. Report your anticipated percentage savings with PV.

Use the data:

1. Square footage of entire house = 2450 sq. ft.
2. Square footage of first floor = 1380 sq. ft.

## Q8. Fill in the gaps below. Let's say you have $18,000 to spend on a PV and get $2,500 in solar rebate . How long will it take (in years) for you to recover the initial investment? (20 points)

```
In [ ]:  area_house = 2450
         area_first_floor = 1380

         XX = dfData_houses['area'].values.reshape(-1, 1)
         YY = dfData_houses['consumption_bill'].values.reshape(-1, 1)

         model_consumption = linear_model.LinearRegression()
         model_consumption.fit(XX, YY)

         print("Predicted power bill = $%.2f" %  model_consumption.predict(np.array([ar

         # Enter code here to compute predicted savings from PV

         del XX, YY
```