# Predicting_bike_sharing_data

September 14, 2020

## 1 Your first neural network

In this project, you'll build your first neural network and use it to predict daily bike rental ridership. We've provided some of the code, but left the implementation of the neural network up to you (for the most part). After you've submitted this project, feel free to explore the data and the model more.

```
[1]: %matplotlib inline
%load_ext autoreload
%autoreload 2
%config InlineBackend.figure_format = 'retina'

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The text.latex.preview rcparam was deprecated in Matplotlib 3.3 and will be
removed two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The mathtext.fallback_to_cm rcparam was deprecated in Matplotlib 3.3 and will be
removed two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle: Support for setting the
'mathtext.fallback_to_cm' rcParam is deprecated since 3.3 and will be removed
two minor releases later; use 'mathtext.fallback : 'cm' instead.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The validate_bool_maybe_none function was deprecated in Matplotlib 3.3 and will
be removed two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The savefig.jpeg_quality rcparam was deprecated in Matplotlib 3.3 and will be
removed two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:

The keymap.all_axes rcparam was deprecated in Matplotlib 3.3 and will be removed
two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The animation.avconv_path rcparam was deprecated in Matplotlib 3.3 and will be
removed two minor releases later.
In /home/user21st/.local/lib/python3.6/site-packages/matplotlib/mpl-
data/stylelib/_classic_test.mplstyle:
The animation.avconv_args rcparam was deprecated in Matplotlib 3.3 and will be
removed two minor releases later.

## 1.1 Load and prepare the data

A critical step in working with neural networks is preparing the data correctly. Variables on different
scales make it difficult for the network to efficiently learn the correct weights. Below, we've written
the code to load and prepare the data. You'll learn more about this soon!

```
[2]: data_path = 'Bike-Sharing-Dataset/hour.csv'

rides = pd.read_csv(data_path)
```

```
[3]: rides.head()
```

```
[3]:    instant      dteday  season  yr  mnth  hr  holiday  weekday  workingday  \
    0        1  2011-01-01       1   0     1   0        0        6           0
    1        2  2011-01-01       1   0     1   1        0        6           0
    2        3  2011-01-01       1   0     1   2        0        6           0
    3        4  2011-01-01       1   0     1   3        0        6           0
    4        5  2011-01-01       1   0     1   4        0        6           0

       weathersit  temp   atemp   hum  windspeed  casual  registered  cnt
    0           1  0.24  0.2879  0.81        0.0       3          13   16
    1           1  0.22  0.2727  0.80        0.0       8          32   40
    2           1  0.22  0.2727  0.80        0.0       5          27   32
    3           1  0.24  0.2879  0.75        0.0       3          10   13
    4           1  0.24  0.2879  0.75        0.0       0           1    1
```

## 1.2 Checking out the data

This dataset has the number of riders for each hour of each day from January 1 2011 to December
31 2012. The number of riders is split between casual and registered, summed up in the `cnt` column.
You can see the first few rows of the data above.

Below is a plot showing the number of bike riders over the first 10 days or so in the data set. (Some
days don't have exactly 24 entries in the data set, so it's not exactly 10 days.) You can see the
hourly rentals here. This data is pretty complicated! The weekends have lower over all ridership
and there are spikes when people are biking to and from work during the week. Looking at the

data above, we also have information about temperature, humidity, and windspeed, all of these likely affecting the number of riders. You'll be trying to capture all this with your model.

```
[4]: rides.describe()
```

```
[4]:            instant        season           yr          mnth            hr  \
       count  17379.0000  17379.000000  17379.000000  17379.000000  17379.000000
       mean    8690.0000      2.501640      0.502561      6.537775     11.546752
       std     5017.0295      1.106918      0.500008      3.438776      6.914405
       min        1.0000      1.000000      0.000000      1.000000      0.000000
       25%     4345.5000      2.000000      0.000000      4.000000      6.000000
       50%     8690.0000      3.000000      1.000000      7.000000     12.000000
       75%    13034.5000      3.000000      1.000000     10.000000     18.000000
       max    17379.0000      4.000000      1.000000     12.000000     23.000000

                  holiday       weekday    workingday     weathersit          temp  \
       count  17379.000000  17379.000000  17379.000000  17379.000000  17379.000000
       mean       0.028770      3.003683      0.682721      1.425283      0.496987
       std        0.167165      2.005771      0.465431      0.639357      0.192556
       min        0.000000      0.000000      0.000000      1.000000      0.020000
       25%        0.000000      1.000000      0.000000      1.000000      0.340000
       50%        0.000000      3.000000      1.000000      1.000000      0.500000
       75%        0.000000      5.000000      1.000000      2.000000      0.660000
       max        1.000000      6.000000      1.000000      4.000000      1.000000

                    atemp           hum     windspeed        casual     registered  \
       count  17379.000000  17379.000000  17379.000000  17379.000000  17379.000000
       mean       0.475775      0.627229      0.190098     35.676218    153.786869
       std        0.171850      0.192930      0.122340     49.305030    151.357286
       min        0.000000      0.000000      0.000000      0.000000      0.000000
       25%        0.333300      0.480000      0.104500      4.000000     34.000000
       50%        0.484800      0.630000      0.194000     17.000000    115.000000
       75%        0.621200      0.780000      0.253700     48.000000    220.000000
       max        1.000000      1.000000      0.850700    367.000000    886.000000

                      cnt
       count  17379.000000
       mean     189.463088
       std      181.387599
       min        1.000000
       25%       40.000000
       50%      142.000000
       75%      281.000000
       max      977.000000
```

```
[5]: rides.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 17379 entries, 0 to 17378
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   instant     17379 non-null  int64
 1   dteday      17379 non-null  object
 2   season      17379 non-null  int64
 3   yr          17379 non-null  int64
 4   mnth        17379 non-null  int64
 5   hr          17379 non-null  int64
 6   holiday     17379 non-null  int64
 7   weekday     17379 non-null  int64
 8   workingday  17379 non-null  int64
 9   weathersit  17379 non-null  int64
 10  temp        17379 non-null  float64
 11  atemp       17379 non-null  float64
 12  hum         17379 non-null  float64
 13  windspeed   17379 non-null  float64
 14  casual      17379 non-null  int64
 15  registered  17379 non-null  int64
 16  cnt         17379 non-null  int64
dtypes: float64(4), int64(12), object(1)
memory usage: 2.3+ MB
```
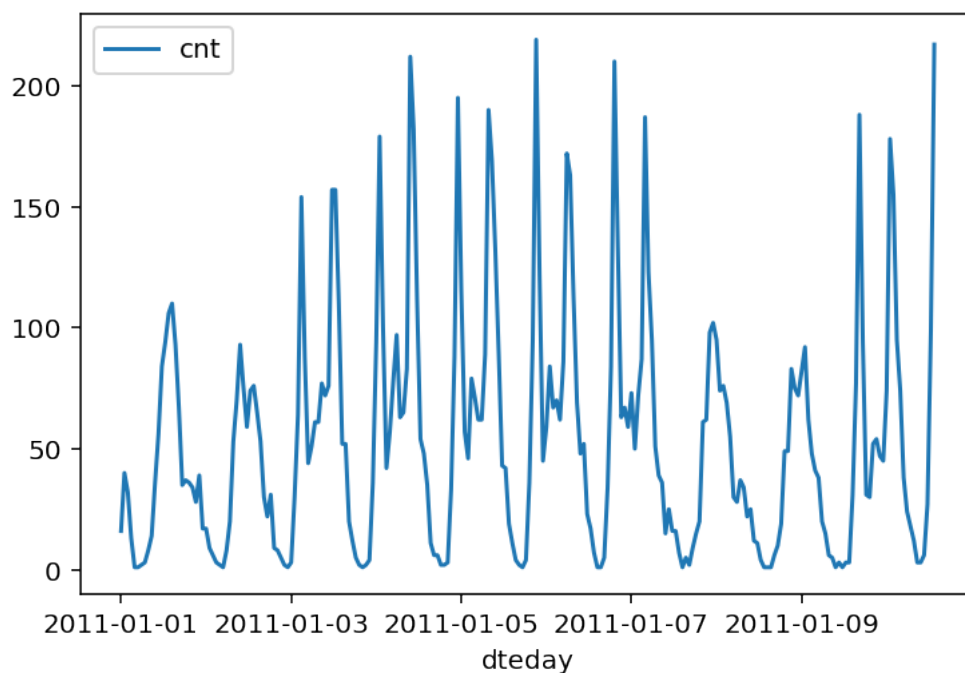
[6]: `rides[:24*10].plot(x='dteday', y='cnt')`

[6]: `<AxesSubplot:xlabel='dteday'>`

### 1.2.1 Dummy variables

Here we have some categorical variables like season, weather, month. To include these in our model, we'll need to make binary dummy variables. This is simple to do with Pandas thanks to `get_dummies()`.

```python
dummy_fields = ['season', 'weathersit', 'mnth', 'hr', 'weekday']
for each in dummy_fields:
    dummies = pd.get_dummies(rides[each], prefix=each, drop_first=False)
    rides = pd.concat([rides, dummies], axis=1)

fields_to_drop = ['instant', 'dteday', 'season', 'weathersit',
                  'weekday', 'atemp', 'mnth', 'workingday', 'hr']
data = rides.drop(fields_to_drop, axis=1)
data.head()
```

[7]:

|   | yr | holiday | temp | hum | windspeed | casual | registered | cnt | season_1 \ |
|---|----|---------|------|-----|-----------|--------|------------|-----|------------|
| 0 | 0 | 0 | 0.24 | 0.81 | 0.0 | 3 | 13 | 16 | 1 |
| 1 | 0 | 0 | 0.22 | 0.80 | 0.0 | 8 | 32 | 40 | 1 |
| 2 | 0 | 0 | 0.22 | 0.80 | 0.0 | 5 | 27 | 32 | 1 |
| 3 | 0 | 0 | 0.24 | 0.75 | 0.0 | 3 | 10 | 13 | 1 |
| 4 | 0 | 0 | 0.24 | 0.75 | 0.0 | 0 | 1 | 1 | 1 |

|   | season_2 | ... | hr_21 | hr_22 | hr_23 | weekday_0 | weekday_1 | weekday_2 \ |
|---|----------|-----|-------|-------|-------|-----------|-----------|-------------|
| 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |

|   | weekday_3 | weekday_4 | weekday_5 | weekday_6 |
|---|-----------|-----------|-----------|-----------|
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 0 | 1 |
| 2 | 0 | 0 | 0 | 1 |
| 3 | 0 | 0 | 0 | 1 |
| 4 | 0 | 0 | 0 | 1 |

[5 rows x 59 columns]

### 1.2.2 Scaling target variables

To make training the network easier, we'll standardize each of the continuous variables. That is, we'll shift and scale the variables such that they have zero mean and a standard deviation of 1.

The scaling factors are saved so we can go backwards when we use the network for predictions.

```
[8]: quant_features = ['casual', 'registered', 'cnt', 'temp', 'hum', 'windspeed']
     # Store scalings in a dictionary so we can convert back later
     scaled_features = {}
     for each in quant_features:
         mean, std = data[each].mean(), data[each].std()
         scaled_features[each] = [mean, std]
         data.loc[:, each] = (data[each] - mean)/std
```

### 1.2.3 Splitting the data into training, testing, and validation sets

We'll save the data for the last approximately 21 days to use as a test set after we've trained the network. We'll use this set to make predictions and compare them with the actual number of riders.

```
[9]: # Save data for approximately the last 21 days
     test_data = data[-21*24:]

     # Now remove the test data from the data set
     data = data[:-21*24]

     # Separate the data into features and targets
     target_fields = ['cnt', 'casual', 'registered']
     features, targets = data.drop(target_fields, axis=1), data[target_fields]
     test_features, test_targets = test_data.drop(target_fields, axis=1),␣
      ↪test_data[target_fields]
```

We'll split the data into two sets, one for training and one for validating as the network is being trained. Since this is time series data, we'll train on historical data, then try to predict on future data (the validation set).

```
[10]: # Hold out the last 60 days or so of the remaining data as a validation set
      train_features, train_targets = features[:-60*24], targets[:-60*24]
      val_features, val_targets = features[-60*24:], targets[-60*24:]
```

## 1.3 Time to build the network

Below you'll build your network. We've built out the structure. You'll implement both the forward pass and backwards pass through the network. You'll also set the hyperparameters: the learning rate, the number of hidden units, and the number of training passes.

The network has two layers, a hidden layer and an output layer. The hidden layer will use the sigmoid function for activations. The output layer has only one node and is used for the regression, the output of the node is the same as the input of the node. That is, the activation function is $f(x) = x$. A function that takes the input signal and generates an output signal, but takes into account the threshold, is called an activation function. We work through each layer of our network

calculating the outputs for each neuron. All of the outputs from one layer become inputs to the neurons on the next layer. This process is called *forward propagation*.

We use the weights to propagate signals forward from the input to the output layers in a neural network. We use the weights to also propagate error backwards from the output back into the network to update our weights. This is called *backpropagation*.

> **Hint:** You'll need the derivative of the output activation function ($f(x) = x$) for the backpropagation implementation. If you aren't familiar with calculus, this function is equivalent to the equation $y = x$. What is the slope of that equation? That is the derivative of $f(x)$.

Below, you have these tasks: 1. Implement the sigmoid function to use as the activation function. Set `self.activation_function` in `__init__` to your sigmoid function. 2. Implement the forward pass in the `train` method. 3. Implement the backpropagation algorithm in the `train` method, including calculating the output error. 4. Implement the forward pass in the `run` method.

```
[11]: #############
      # In the neural_network.py file, fill out the TODO sections as specified
      #############


      from neural_network import NeuralNetwork
```

```
[12]: def MSE(y, Y):
          return np.mean((y-Y)**2)
```

## 1.4 Unit tests

Run these unit tests to check the correctness of your network implementation. This will help you be sure your network was implemented correctly befor you starting trying to train it. These tests must all be successful to pass the project.

```
[13]: import unittest


      inputs = np.array([[0.5, -0.2, 0.1]])
      targets = np.array([[0.4]])
      test_w_i_h = np.array([[0.1, -0.2],
                             [0.4, 0.5],
                             [-0.3, 0.2]])
      test_w_h_o = np.array([[0.3],
                             [-0.1]])


      class TestMethods(unittest.TestCase):

          ##########
          # Unit tests for data loading
          ##########
```

```python
    def test_data_path(self):
        # Test that file path to dataset has been unaltered
        self.assertTrue(data_path.lower() == 'bike-sharing-dataset/hour.csv')

    def test_data_loaded(self):
        # Test that data frame loaded
        self.assertTrue(isinstance(rides, pd.DataFrame))

    ##########
    # Unit tests for network functionality
    ##########

    def test_activation(self):
        network = NeuralNetwork(3, 2, 1, 0.5)
        # Test that the activation function is a sigmoid
        self.assertTrue(np.all(network.activation_function(0.5) == 1/(1+np.
 exp(-0.5))))

    def test_train(self):
        # Test that weights are updated correctly on training
        network = NeuralNetwork(3, 2, 1, 0.5)
        network.weights_input_to_hidden = test_w_i_h.copy()
        network.weights_hidden_to_output = test_w_h_o.copy()

        network.train(inputs, targets)
        self.assertTrue(np.allclose(network.weights_hidden_to_output,
                                    np.array([[ 0.37275328],
                                              [-0.03172939]])))
        self.assertTrue(np.allclose(network.weights_input_to_hidden,
                                    np.array([[ 0.10562014, -0.20185996],
                                              [0.39775194, 0.50074398],
                                              [-0.29887597, 0.19962801]])))

    def test_run(self):
        # Test correctness of run method
        network = NeuralNetwork(3, 2, 1, 0.5)
        network.weights_input_to_hidden = test_w_i_h.copy()
        network.weights_hidden_to_output = test_w_h_o.copy()

        self.assertTrue(np.allclose(network.run(inputs), 0.09998924))

suite = unittest.TestLoader().loadTestsFromModule(TestMethods())
unittest.TextTestRunner().run(suite)
```

…
```
----------------------------------------------------------------------
Ran 5 tests in 0.003s
```

```
OK
```

```
[13]: <unittest.runner.TextTestResult run=5 errors=0 failures=0>
```

## 1.5 Training the network

Here you'll set the hyperparameters for the network. The strategy here is to find hyperparameters such that the error on the training set is low, but you're not overfitting to the data. If you train the network too long or have too many hidden nodes, it can become overly specific to the training set and will fail to generalize to the validation set. That is, the loss on the validation set will start increasing as the training set loss drops.

You'll also be using a method know as Stochastic Gradient Descent (SGD) to train the network. The idea is that for each training pass, you grab a random sample of the data instead of using the whole data set. You use many more training passes than with normal gradient descent, but each pass is much faster. This ends up training the network more efficiently. You'll learn more about SGD later.

### 1.5.1 Choose the number of iterations

This is the number of batches of samples from the training data we'll use to train the network. The more iterations you use, the better the model will fit the data. However, this process can have sharply diminishing returns and can waste computational resources if you use too many iterations. You want to find a number here where the network has a low training loss, and the validation loss is at a minimum. The ideal number of iterations would be a level that stops shortly after the validation loss is no longer decreasing.

### 1.5.2 Choose the learning rate

This scales the size of weight updates. If this is too big, the weights tend to explode and the network fails to fit the data. Normally a good choice to start at is 0.1; however, if you effectively divide the learning rate by n_records, try starting out with a learning rate of 1. In either case, if the network has problems fitting the data, try reducing the learning rate. Note that the lower the learning rate, the smaller the steps are in the weight updates and the longer it takes for the neural network to converge.

### 1.5.3 Choose the number of hidden nodes

In a model where all the weights are optimized, the more hidden nodes you have, the more accurate the predictions of the model will be. (A fully optimized model could have weights of zero, after all.) However, the more hidden nodes you have, the harder it will be to optimize the weights of the model, and the more likely it will be that suboptimal weights will lead to overfitting. With overfitting, the model will memorize the training data instead of learning the true pattern, and won't generalize well to unseen data.

Try a few different numbers and see how it affects the performance. You can look at the losses dictionary for a metric of the network performance. If the number of hidden units is too low, then the model won't have enough space to learn and if it is too high there are too many options for the direction that the learning can take. The trick here is to find the right balance in number of hidden units you choose. You'll generally find that the best number of hidden nodes to use ends up being between the number of input and output nodes.

```python
import sys


####################
### Set the hyperparameters in you neural_network.py file ###
####################

from neural_network import iterations, learning_rate, hidden_nodes, output_nodes


N_i = train_features.shape[1]
network = NeuralNetwork(N_i, hidden_nodes, output_nodes, learning_rate)

losses = {'train':[], 'validation':[]}
for ii in range(iterations):
    # Go through a random batch of 128 records from the training data set
    batch = np.random.choice(train_features.index, size=128)
    X, y = train_features.iloc[batch].values, train_targets.iloc[batch]['cnt']

    network.train(X, y)

    # Printing out the training progress
    train_loss = MSE(network.run(train_features).T, train_targets['cnt'].values)
    val_loss = MSE(network.run(val_features).T, val_targets['cnt'].values)
    sys.stdout.write("\rProgress: {:2.1f}".format(100 * ii/float(iterations)) \
                     + "% ... Training loss: " + str(train_loss)[:5] \
                     + " ... Validation loss: " + str(val_loss)[:5])
    sys.stdout.flush()

    losses['train'].append(train_loss)
    losses['validation'].append(val_loss)
```
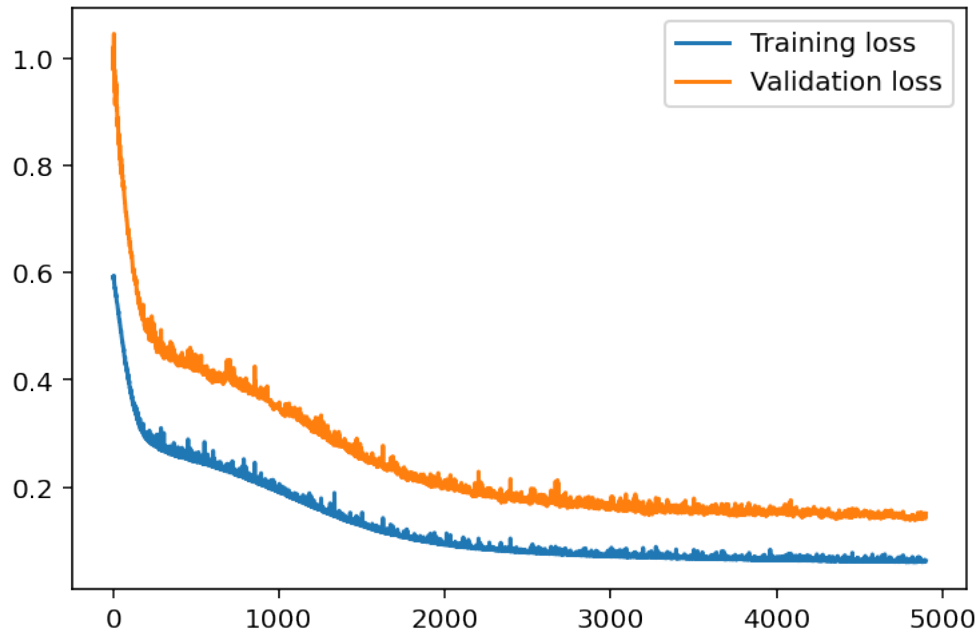
Progress: 100.0% … Training loss: 0.062 … Validation loss: 0.150

```python
plt.plot(losses['train'][100:], label='Training loss')
plt.plot(losses['validation'][100:], label='Validation loss')
plt.legend()
_ = plt.ylim()
```
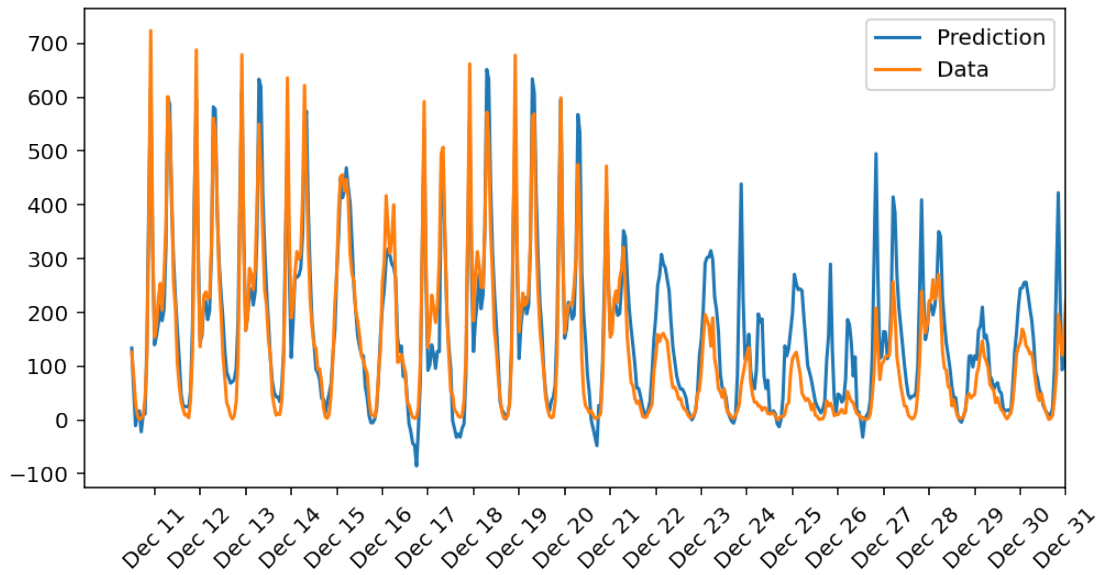
## 1.6 Check out your predictions

Here, use the test data to view how well your network is modeling the data. If something is completely wrong here, make sure each step in your network is implemented correctly.

```
[16]: fig, ax = plt.subplots(figsize=(8,4))

mean, std = scaled_features['cnt']
predictions = network.run(test_features).T*std + mean
ax.plot(predictions[0], label='Prediction')
ax.plot((test_targets['cnt']*std + mean).values, label='Data')
ax.set_xlim(right=len(predictions))
ax.legend()

dates = pd.to_datetime(rides.iloc[test_data.index]['dteday'])
dates = dates.apply(lambda d: d.strftime('%b %d'))
ax.set_xticks(np.arange(len(dates))[12::24])
_ = ax.set_xticklabels(dates[12::24], rotation=45)
```

## 1.7 OPTIONAL: Thinking about your results(this question will not be evaluated in the rubric).

Answer these questions about your results. How well does the model predict the data? Where does it fail? Why does it fail where it does?

The model predictions seem to follow closely with test data labels but it doesnot vary with the same magnitude as the labels.