

Predikce pomoci ANN

```
In [1]: import numpy as np
from neural_net.neuralnet import NeuralNet
import neural_net.activation_functions as af
import pandas as pd
from matplotlib import pyplot as plt
import matplotlib.pyplot as plt
import random
import csv
from pylab import rcParams
rcParams['figure.figsize'] = 16, 6 # resize figures
plt.style.use('bmh')
```

Úvod

Tento projekt se pokouší predikovat vývoj veličiny avg temps na základě 5 proměnných, ze kterých není veličina přímo odvozená:

- Precipitation
- Insolation
- Evaporation
- AvgHumidity
- Windspeed

Preprocessing dat

```
In [2]: with open('resources/cruzeirodosul2018daily.csv', 'r') as f:
    avg_temps = csv.DictReader(f, delimiter=',')
    avg_temps = np.array([i['AvgTemp'] for i in reader])
    avg_temps = np.array([np.nan if i == "" else i for i in avg_temps], dtype='float64') # replace missing val
with open('resources/cruzeirodosul2018daily.csv', 'r') as f:
    reader = csv.DictReader(f, delimiter=',')
    dates = np.array([i['Date'] for i in reader], dtype='int64') / 100.000 # squish the range of dates
with open('resources/cruzeirodosul2018daily.csv', 'r') as f:
    reader = csv.DictReader(f, delimiter=',')
    inputs = np.array([i['Precipitation'], i['Insolation'], i['Evaporation'], i['AvgHumidity'], i['WindSpeed']]
    inputs = np.where(inputs == '', np.nan, inputs)
    inputs = inputs.astype('float64')

dates[0:5], avg_temps[0:5], inputs
Out[2]: (array([[0.40179, 0.4018, 0.40181, 0.40182, 0.40183],
        array([[27.414, 27.7, 28.68, 26.4, 24.98],
        avg_temps[0:5], (array([[0.0, 0.5, 1.7, 79.25, 1.66666],
        [0., 0.7, 3.5, 75.25, 1.],
        [0., 0.3, 4.1, 73., 0.],
        ...,
        [5.8, 0.8, nan, 92., 1.82888],
        [0., 0.2, nan, 93., 0.],
        [0., 7.1, nan, 97., 0. ]]))

Grafické znázornění časové řady
```

```
In [3]: plt.plot(dates, avg_temps)

Out[3]: <matplotlib.lines.Line2D at 0x7f80f832400>
```

Doplnění chybějících hodnot

Ověření předpokladu, že se jedná o časovou řadu

```
In [4]: np.all(np.diff(dates) >= 0)

Out[4]: True

Lineární interpolace chybějících dat k výstupů, vypočtení popisné statistiky sřídění hodnoty a rozptyl (pro 0 stupňů volnosti)

In [5]: avg_temps_interp = NeuralNet.linear_interpolation(avg_temps)
random.shuffle(avg_temps_interp, np.var(avg_temps_interp))

Out[5]: (25.867502343817804, 1.9603880733839343)

In [6]: len(inputs), len(avg_temps_interp)

Out[6]: (2134, 2134)

Aplikování normalizace, detrendování diferenciál ř. řady. Chybějící data doplněná lineární interpolací

In [7]: norm_avg_temps = NeuralNet.normalize(avg_temps_interp)
diff_norm_avg_temps = np.diff(norm_avg_temps)

inputs_interp = NeuralNet.linear_interpolation(inputs)
norm_inputs = NeuralNet.normalize(inputs_interp)
diff_norm_inputs = np.diff(norm_inputs, axis=0)

norm_inputs, diff_norm_inputs

Out[7]: (array([[0., 0.887236, 0.01561965, 0.72773186, 0.01530456],
        [0., 0.07888881, 0.03213958, 0.69180892, 0.00918274],
        [0., 0.07621671, 0.03764922, 0.67833976, 0.],
        ...,
        [0.95325987, 0.00880888, 0.00944793, 0.84811175, 0.00944793],
        [0., 0.02835055, 0.01838555, 0.00472397, 0.04591368, 0.00472397],
        [0., 0.06519743, 0., 0.8972544, 0.],
        [0., 0.05122991],
        array([[0.0228476, -0.01469238, 0.00734619, 0.05059095, 0.01836547],
        [0.00734619, 0.01652893, -0.03673895, -0.00612183],
        [0.0228476, -0.00867389, 0.00560064, -0.02866116, 0.00918274],
        [0.04224859, -0.06336888, -0.00891827, 0.07885326, 0.00918274],
        ...,
        [-0.05325987, 0.05785124, 0.00472397, 0.03673895, 0.00472397],
        [-0.05325987, -0.07897153, -0.00944793, 0.00918274, -0.00944793],
        [0., 0.06336888, 0., 0.03673895, 0. ]]))

In [8]: plt.plot(dates[1:], diff_norm_avg_temps)

Out[8]: <matplotlib.lines.Line2D at 0x7f80d7ddf78>
```

Úprava dat pro neuronovou síť

množina byla rozdělena na trénovací a validační v poměru 75%:25%

```
In [9]: dataset = [(x, np.array(y))] for x, y in zip(diff_norm_inputs, diff_norm_avg_temps)]
random.shuffle(dataset)
dataset[0:5]

Out[9]: ([array([0.4224859, -0.0174472, -0.00642792, 0.07116621, 0.],
        array([-0.01813298])),
        (array([-0.11753983, -0.07162534, 0.00472397, -0.04591368, 0.00472397]),
        [0., 0.08738891]),
        (array([-0.1423242, 0.00183655, 0.00734619, 0.05059095, 0.01836547]),
        [0., 0.06519743, 0., 0.8972544, 0.]),
        (array([-0.03122991]),
        [0., 0.0228476, -0.01469238, 0.00734619, 0.05059095, 0.01836547]),
        (array([-0.05104097]),
        [0.0228476, -0.01469238, 0., -0.03673895, 0.]),
        (array([0.00375796]))])

In [10]: training_data, test_data = NeuralNet.ratio_list_split(dataset, 0.75)
len(training_data), len(test_data), len(dataset), training_data[0]

Out[10]: (1599,
534,
2133,
(array([0.4224859, -0.0174472, -0.00642792, 0.07116621, 0.],
        array([-0.01813298]))))
```

Vypracování

Připrava množiny hyperparametrů

```
In [11]: neurons = [3, 5, 10, 50]
epochs = [True]
learning_rates = [0.5, 0.1, 0.05, 0.005]
max_epochs = 200
target_test_mse = 0.001

In [12]: params = itertools.product(neurons, epochs, learning_rates)
stats = []
for p in params:
    test_mse_list = []
    train_mse_list = []
    nn = NeuralNet([len(dataset[0][0]), p[0], len(dataset[0][1])], a_functions=[af.tanh, af.tanh], a_function=print)
    for e in nn.gradient_descent_testdata(
        labeled_training_dataset=training_data,
        labeled_test_dataset=test_data,
        no_epochs=p[1],
        mini_batch_size=128,
        learning_rate=p[2]
    ):
        print('epoch: ', e['epoch'], end='\n')
        test_mse_list.append(e['test_mse'])
        train_mse_list.append(e['train_mse'])
        if e['test_mse'] < target_test_mse:
            result = {
                'success': True,
                'neurons input layer': nn.layer_sizes[0],
                'neurons hidden layer': p[0],
                'epochs': e['epoch'] + 1,
                'learning rate': p[2],
                'test_mse list': test_mse_list,
                'train_mse list': train_mse_list,
            }
            stats.append(result)
            print()
            print('TARGET HIT epoch', e['epoch'], '\n test_mse', e['test_mse'],
            print('\n train_mse', e['train_mse'])
            print()
            break
        if e['epoch'] > max_epochs - 1: # last epoch
            result = {
                'success': False,
                'neurons input layer': nn.layer_sizes[0],
                'neurons hidden layer': p[0],
                'epochs': e['epoch'] + 1,
                'learning rate': p[2],
                'test_mse list': test_mse_list,
                'train_mse list': train_mse_list,
            }
            stats.append(result)
            print()
            print('TARGET NOT HIT epoch', e['epoch'], '\n test_mse', e['test_mse'],
            print('\n train_mse', e['train_mse'])
            print()
            break
stats_json = json.dumps(stats)
with open('resources/stats_avgtemp.json', 'w') as f:
    f.write(stats_json)

(3, True, 0.5)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.001208097790706013
train_mse 0.362395404406406e-06

(3, True, 0.1)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.002071730270952134
train_mse 1.6075633671463267e-05

(3, True, 0.05)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.001995691842135233
train_mse 1.5771827640289893e-05

(3, True, 0.005)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.004111418215085452
train_mse 3.8092954309637197e-05

(5, True, 0.5)
epoch: 200182
TARGET NOT HIT epoch 200
test_mse 0.001184012863767278
train_mse 1.2127280620917204e-05

(5, True, 0.1)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.00151523929895873
train_mse 1.237080057576866e-05

(5, True, 0.05)
epoch: 200118
TARGET NOT HIT epoch 200
test_mse 0.001766838080845378
train_mse 1.4468875444265574e-05

(5, True, 0.005)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.0037766321628151943
train_mse 3.80287576447289e-05

(10, True, 0.5)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.00181523929895873
train_mse 0.000925697540e-06

(10, True, 0.1)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.001294821731240695
train_mse 1.044363325126191e-05

(10, True, 0.05)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.00197224828561035
train_mse 1.6439216180609204e-05

(10, True, 0.005)
epoch: 2003
TARGET NOT HIT epoch 200
test_mse 0.005722682806281198
train_mse 7.744383084404097e-05

(50, True, 0.5)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.001193274356787623
train_mse 7.40148288684450e-06

(50, True, 0.1)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.001918092471827400
train_mse 1.3740396755756245e-05

(50, True, 0.05)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.0413256808471941
train_mse 0.0003674478165181378

(50, True, 0.005)
epoch: 200
TARGET NOT HIT epoch 200
test_mse 0.021456680630832176
train_mse 0.000153425634078429
```

Citlivostní analýza

tabulka hyperparametrů

```
In [13]: df = pd.read_json('resources/stats_avgtemp.json')
df[['neurons input layer', 'neurons hidden layer', 'epochs', 'learning rate']]

Out[13]:
```

neurons input layer	neurons hidden layer	epochs	learning rate
5	3	201	0.500
1	5	3	0.100
2	5	3	0.050
3	5	3	0.005
4	5	5	0.500
5	5	5	0.100
6	5	5	0.050
7	5	5	0.005
8	5	10	0.500
9	5	10	0.100
10	5	10	0.050
11	5	10	0.005
12	5	50	0.500
13	5	50	0.100
14	5	50	0.050
15	5	50	0.005

```
In [14]: # fig.set_xlim((0, 200))
# fig.set_ylim((0, 0.05))
for p in range(15):
    plt.plot(df['train_mse_list'][p], label='train mse')
    plt.plot(df['test_mse_list'][p], label='test mse')
    plt.title(f'input layer neurons: {p} | hidden layer neurons: {}, eta: {}'.format(df['neurons input layer'],
    plt.xlabel('epoch')
    plt.ylabel('mse')
    plt.legend()
    ax, fig = plt.subplots()
    plt.plot(df['train_mse_list'][p], label='train mse')
    plt.plot(df['test_mse_list'][p], label='test mse')
    plt.title(f'input layer neurons: {p} | hidden layer neurons: {}, eta: {}'.format(df['neurons input layer'],
    plt.xlabel('epoch')
    plt.ylabel('mse')
    plt.legend()

Out[14]: <matplotlib.legend.Legend at 0x7f80a427b310>
```

Po vizuální inspekci průběhu vývoje test mse byly pro predikci vybrány hyperparametry:

- počet neuronů vstupů: 4
- počet neuronů skrytá vrstva: 50
- learning rate 0.05

```
In [15]: df = pd.read_json('resources/stats_avgtemp_bestapprox_targetlower.json')
ax, fig = plt.subplots()
# fig.set_xlim((0, 0.02))
plt.plot(df['train_mse_list'][0], label='train mse')
plt.plot(df['test_mse_list'][0], label='test mse')
plt.title('input layer neurons: 4 | hidden layer neurons: 50, eta: 0.05')
plt.xlabel('epoch')
plt.ylabel('mse')
plt.legend()

Out[15]: <matplotlib.legend.Legend at 0x7f80d069420>
```

```
In [19]: prediction = nn.feed_forward(x)['x'] for x in diff_norm_inputs]
prediction = np.array([x[0] for x in prediction])
prediction, len(prediction)

Out[19]: (array([-0.0820993, -0.06500859, -0.12357382, ..., -0.04517385,
        -0.08289188, -0.04136453]),
2133)

In [20]: prediction

Out[20]: array([-0.0820993, -0.06500859, -0.12357382, ..., -0.04517385,
        -0.08289188, -0.04136453])
```

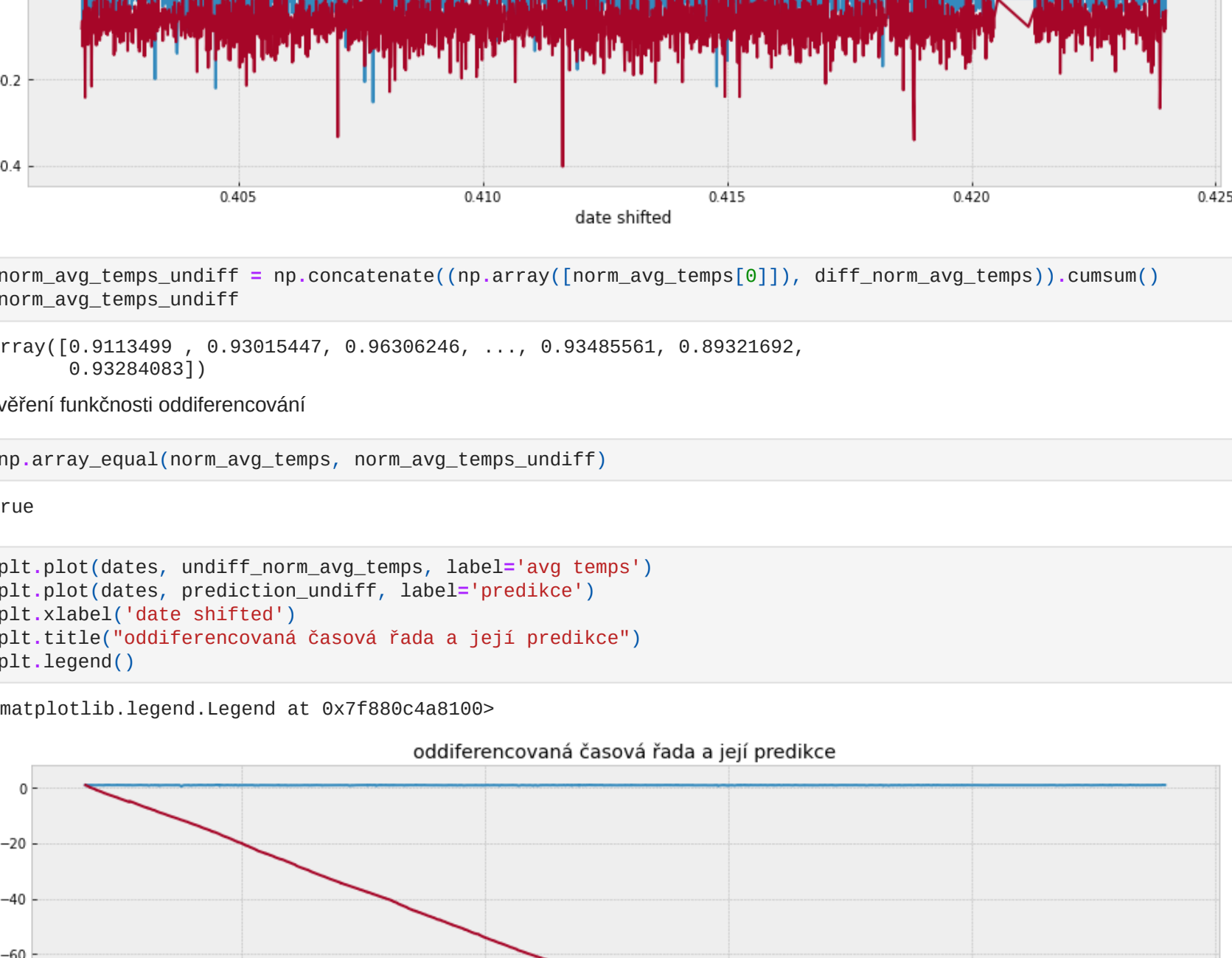


```
In [25]: prediction_undiff = np.concatenate((np.array([norm_avg_temps[0]]), prediction)).cumsum()

Out[25]: array([ 0.0113499, 0.8293486, 0.76433291, ..., -140.77739914,
-140.86629103, -140.90165556])

In [37]: # plt.plot(dates, norm_avg_temps)
plt.plot(dates[1:], diff_norm_avg_temps, label='difference avg temps')
plt.plot(dates[1:], prediction, label='predikce difference avg temps')
plt.xlabel('date shifted')
plt.title("difference časové řady a její predikce")
plt.legend()

Out[37]: <matplotlib.legend.Legend at 0x7f680c1a9460>
```



```
In [31]: norm_avg_temps_undiff = np.concatenate((np.array([norm_avg_temps[0]]), diff_norm_avg_temps)).cumsum()
norm_avg_temps_undiff

Out[31]: array([0.0113499, 0.93915447, 0.96366246, ..., 0.93485561, 0.89321692,
0.93284983])


Ověření funkčnosti oddiferencování

In [32]: np.array_equal(norm_avg_temps, norm_avg_temps_undiff)

Out[32]: True

In [38]: plt.plot(dates, undiff_norm_avg_temps, label='avg temps')
plt.plot(dates, prediction_undiff, label='predikce')
plt.xlabel('date shifted')
plt.title('oddiferencovaná časová řada a její predikce')
plt.legend()

Out[38]: <matplotlib.legend.Legend at 0x7f680c4a8100>
```



Závěr

Po vizuální inspekci difference v predikci byl nalezen posun na ose y do záporných hodnot, z čehož vyplývá velký klesající trend při predikci celé časové řady.