

E14 BP Algorithm (C++/Python)

17341015 Hongzheng Chen

December 14, 2019

Contents

1	Horse Colic Data Set	2
2	Reference Materials	2
3	Tasks	5
4	Codes and Results	6

1 Horse Colic Data Set

The description of the horse colic data set (<http://archive.ics.uci.edu/ml/datasets/Horse+Colic>) is as follows:

Data Set Characteristics:	Multivariate	Number of Instances:	368	Area:	Life
Attribute Characteristics:	Categorical, Integer, Real	Number of Attributes:	27	Date Donated	1989-08-06
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	108569

We aim at trying to predict if a horse with colic will live or die.

Note that we should deal with missing values in the data! Here are some options:

- Use the features mean value from all the available data.
- Fill in the unknown with a special value like -1.
- Ignore the instance.
- Use a mean value from similar items.
- Use another machine learning algorithm to predict the value.

2 Reference Materials

1. Stanford: **CS231n: Convolutional Neural Networks for Visual Recognition** by Fei-Fei Li, etc.
 - Course website: <http://cs231n.stanford.edu/2017/syllabus.html>
 - Video website: https://www.bilibili.com/video/av17204303/?p=9&tdsourcetag=s_pctim_aiomsg
2. **Machine Learning** by Hung-yi Lee
 - Course website: <http://speech.ee.ntu.edu.tw/~tlkagk/index.html>
 - Video website: <https://www.bilibili.com/video/av9770302/from=search>
3. A Simple neural network code template

```
1 # -*- coding: utf-8 -*-
2 import random
3 import math
4
5 # Shorthand:
6 # "pd_" as a variable prefix means "partial derivative"
7 # "d_" as a variable prefix means "derivative"
8 # "_wrt_" is shorthand for "with respect to"
9 # "w_ho" and "w_ih" are the index of weights from hidden to output layer neurons and
  ↳ input to hidden layer neurons respectively
10
11 class NeuralNetwork:
12     LEARNING_RATE = 0.5
13     def __init__(self, num_inputs, num_hidden, num_outputs, hidden_layer_weights =
  ↳ None, hidden_layer_bias = None, output_layer_weights = None,
  ↳ output_layer_bias = None):
14         #Your Code Here
15
16     def init_weights_from_inputs_to_hidden_layer_neurons(self, hidden_layer_weights):
17         #Your Code Here
18
19     def init_weights_from_hidden_layer_neurons_to_output_layer_neurons(self,
  ↳ output_layer_weights):
```

```

20 #Your Code Here
21
22 def inspect(self):
23     print('-----')
24     print('* Inputs: {}'.format(self.num_inputs))
25     print('-----')
26     print('Hidden Layer')
27     self.hidden_layer.inspect()
28     print('-----')
29     print('* Output Layer')
30     self.output_layer.inspect()
31     print('-----')
32
33 def feed_forward(self, inputs):
34     #Your Code Here
35
36 # Uses online learning, ie updating the weights after each training case
37 def train(self, training_inputs, training_outputs):
38     self.feed_forward(training_inputs)
39
40     # 1. Output neuron deltas
41     #Your Code Here
42     #  $E / z$ 
43
44     # 2. Hidden neuron deltas
45     # We need to calculate the derivative of the error with respect to the output
46     #  $\rightarrow$  of each hidden layer neuron
47     #  $dE / dy = E / z * z / y = E / z * w$ 
48     #  $E / z = dE / dy * z / w$ 
49     #Your Code Here
50
51     # 3. Update output neuron weights
52     #  $E / w = E / z * z / w$ 
53     #  $w = * E / w$ 
54     #Your Code Here
55
56     # 4. Update hidden neuron weights
57     #  $E / w = E / z * z / w$ 
58     #  $w = * E / w$ 
59     #Your Code Here
60
61 def calculate_total_error(self, training_sets):
62     #Your Code Here
63     return total_error
64
65 class NeuronLayer:
66     def __init__(self, num_neurons, bias):
67
68         # Every neuron in a layer shares the same bias
69         self.bias = bias if bias else random.random()
70
71         self.neurons = []
72         for i in range(num_neurons):
73             self.neurons.append(Neuron(self.bias))
74
75     def inspect(self):

```

```

75     print('Neurons:', len(self.neurons))
76     for n in range(len(self.neurons)):
77         print(' Neuron', n)
78         for w in range(len(self.neurons[n].weights)):
79             print(' Weight:', self.neurons[n].weights[w])
80         print(' Bias:', self.bias)
81
82     def feed_forward(self, inputs):
83         outputs = []
84         for neuron in self.neurons:
85             outputs.append(neuron.calculate_output(inputs))
86         return outputs
87
88     def get_outputs(self):
89         outputs = []
90         for neuron in self.neurons:
91             outputs.append(neuron.output)
92         return outputs
93
94 class Neuron:
95     def __init__(self, bias):
96         self.bias = bias
97         self.weights = []
98
99     def calculate_output(self, inputs):
100         #Your Code Here
101
102     def calculate_total_net_input(self):
103         #Your Code Here
104
105     # Apply the logistic function to squash the output of the neuron
106     # The result is sometimes referred to as 'net' [2] or 'net' [1]
107     def squash(self, total_net_input):
108         #Your Code Here
109
110     # Determine how much the neuron's total input has to change to move closer to the
111         ↪ expected output
112     #
113     # Now that we have the partial derivative of the error with respect to the output
114         ↪ (E/y) and
115     # the derivative of the output with respect to the total net input (dy/dz) we can
116         ↪ calculate
117     # the partial derivative of the error with respect to the total net input.
118     # This value is also known as the delta () [1]
119     # = E / z = E / y * dy / dz
120     #
121     def calculate_pd_error_wrt_total_net_input(self, target_output):
122         #Your Code Here
123
124     # The error for each neuron is calculated by the Mean Square Error method:
125     def calculate_error(self, target_output):
126         #Your Code Here
127
128     # The partial derivate of the error with respect to actual output then is
129         ↪ calculated by:
130     # = 2 * 0.5 * (target output - actual output) ^ (2 - 1) * -1

```

```

127     # = -(target output - actual output)
128     #
129     # The Wikipedia article on backpropagation [1] simplifies to the following, but
130     # ↪ most other learning material does not [2]
131     # = actual output - target output
132     #
133     # Alternative, you can use (target - output), but then need to add it during
134     # ↪ backpropagation [3]
135     #
136     # Note that the actual output of the output neuron is often written as y and
137     # ↪ target output as t so:
138     # = E / y = -(t - y)
139     def calculate_pd_error_wrt_output(self, target_output):
140     #Your Code Here
141
142     # The total net input into the neuron is squashed using logistic function to
143     # ↪ calculate the neuron's output:
144     # y = 1 / (1 + e^(-z))
145     # Note that where represents the output of the neurons in whatever layer we're
146     # ↪ looking at and represents the layer below it
147     #
148     # The derivative (not partial derivative since there is only one variable) of the
149     # ↪ output then is:
150     # dy / dz = y * (1 - y)
151     def calculate_pd_total_net_input_wrt_input(self):
152     #Your Code Here
153
154     # The total net input is the weighted sum of all the inputs to the neuron and
155     # ↪ their respective weights:
156     # z = net = xw + xw ...
157     #
158     # The partial derivative of the total net input with respect to a given weight
159     # ↪ (with everything else held constant) then is:
160     # z / w = some constant + 1 * xw^(1-0) + some constant ... = x
161     def calculate_pd_total_net_input_wrt_weight(self, index):
162     #Your Code Here
163
164     # An example:
165
166     nn = NeuralNetwork(2, 2, 2, hidden_layer_weights=[0.15, 0.2, 0.25, 0.3],
167     # ↪ hidden_layer_bias=0.35, output_layer_weights=[0.4, 0.45, 0.5, 0.55],
168     # ↪ output_layer_bias=0.6)
169     for i in range(10000):
170     nn.train([0.05, 0.1], [0.01, 0.99])
171     print(i, round(nn.calculate_total_error([[0.05, 0.1], [0.01, 0.99]]), 9))

```

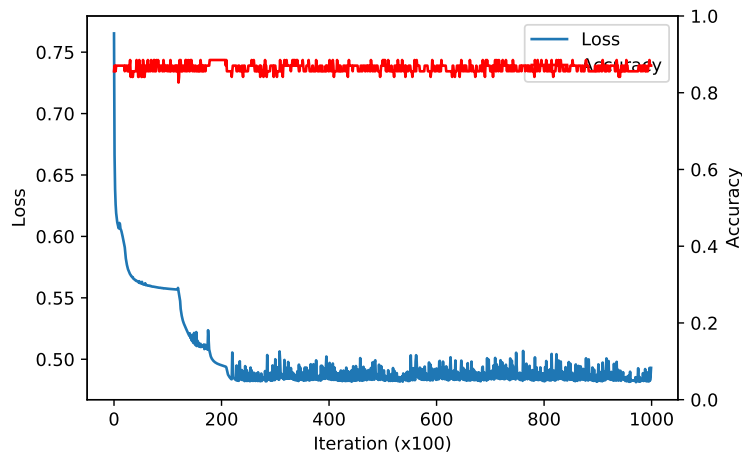
3 Tasks

- Given the training set `horse-colic.data` and the testing set `horse-colic.test`, implement the BP algorithm and establish a neural network to predict if horses with colic will live or die. In addition, you should calculate the accuracy rate.
- Please submit a file named `E14.YourNumber.pdf` and send it to `ai_201901@foxmail.com`

4 Codes and Results

This experiment costs me three full days to finish (finetune the hyperparameters), but I still cannot figure out why my accuracy is so awkward. Sad :(

The following figure gives the training loss and the accuracy (without early stopping).



The training log is shown below (with early stopping), and the best accuracy I can get is 92% accuracy on test set. (Notice the two figures are not in the same training process.)

```
Iter 1000/100000 Loss: 0.6910220450181086 Accuracy: 87.58823529411765%
Iter 1100/100000 Loss: 0.6750272830662208 Accuracy: 87.58823529411765%
Iter 1200/100000 Loss: 0.665812515211816 Accuracy: 87.58823529411765%
Iter 1300/100000 Loss: 0.6597815957555795 Accuracy: 87.58823529411765%
Iter 1400/100000 Loss: 0.6545913532922571 Accuracy: 90.5294117647059%
Iter 1500/100000 Loss: 0.6480593287768979 Accuracy: 92.0%
Iter 1600/100000 Loss: 0.6427443958094696 Accuracy: 92.0%
Iter 1700/100000 Loss: 0.6384380537773697 Accuracy: 92.0%
Iter 1800/100000 Loss: 0.6353141344986356 Accuracy: 92.0%
Iter 1900/100000 Loss: 0.6330266213543354 Accuracy: 92.0%
Iter 2000/100000 Loss: 0.6312329687173412 Accuracy: 92.0%
Iter 2100/100000 Loss: 0.6296858360830062 Accuracy: 92.0%
Iter 2200/100000 Loss: 0.6282317805795642 Accuracy: 92.0%
Iter 2300/100000 Loss: 0.6268011951659986 Accuracy: 92.0%
Early stop at iter 2400/100000 Loss: 0.625693064579745 Accuracy: 90.5294117647059%
Time: 3.130645275115967s
```

Please refer to `nn.py` and `bp.ipynb` (or the generated `bp.py`) for the codes.

Highlight some used techniques:

- The network structure is n_i-8-3 .
- Used Kaiming He's method to initialize the network
- A pytorch-like network class with `forward` method is designed.
- L2-regulization and weight decay are used for training.
- All computation are based on tensor, and only the numpy package is used.
- `np.einsum` is used for accelerating the tensor product in backpropagation.
- Heavy preprocessing methods are used, including one-hot encoding, missing data complement, and useless attributes removal. please refer to `bp.ipynb` file for details.
- Early stopping is used to avoid overfitting, and learning rate decay is used for better convergence.
- Batch SGD is used to accelerate training.
- Checkpoints and logging make training more controllable.

Following gives the code of nn.py.

```
1 import numpy as np
2
3 class FullyConnectedLayer(object):
4     """
5     Linear transformation:  $y = x W^T + b$ 
6     Input: (N, in_features), i.e. a row vector, in this example, N = 1
7     Output: (N, out_features)
8
9     Attributes:
10         Weight: (in_features, out_features)
11         Bias: (out_features)
12
13     Ref:
14     http://cs231n.stanford.edu/vecDerivs.pdf
15     """
16
17     def __init__(self, in_features, out_features, bias=True):
18         self.in_features = in_features
19         self.out_features = out_features
20         """
21         Xavier initialization
22         # https://www.deeplearning.ai/ai-notes/initialization/
23          $W^{[l]} \sim \mathcal{N}(\mu=0, \sigma^2 = \frac{1}{n^{[l-1]}})$ 
24          $b^{[l]} = 0$ 
25
26         Kaiming He initialization
27         # https://medium.com/@shoray.goel/kaiming-he-initialization-a8d9ed0b5899
28         """
29         self.weight = np.random.normal(0, np.sqrt(2/in_features), (out_features, in_features))
30         if bias:
31             self.bias = np.random.rand(out_features)
32         else:
33             self.bias = None
34
35     def forward(self, inputs):
36         """
37         Forward propagation
38         """
39         if type(self.bias) != type(None):
40             return np.dot(inputs, self.weight.T) + self.bias
41         else:
42             return np.dot(inputs, self.weight.T)
43
44     def __call__(self, x):
45         """
46         Syntax sugar for forward method
47         """
48         return self.forward(x)
49
50 class Network(object):
51
52     def __init__(self, in_features, hidden_features, out_features, learning_rate=0.01):
53         """
54         Here three-layer network architecture is used
55         """
```

```

56     The number of neurons in each layer is listed below:
57     in_features -> hidden_features -> out_features
58     """
59     self.fc1 = FullyConnectedLayer(in_features,hidden_features,True)
60     self.fc2 = FullyConnectedLayer(hidden_features,out_features,True)
61     self.learning_rate = learning_rate
62     self.memory = {} # used for store intermediate results
63     self.train_flag = True
64
65     def train(self):
66         """
67         When training, memory is set to remember the intermediate results
68         """
69         self.train_flag = True
70
71     def eval(self):
72         """
73         When inferencing, memory is no need to set
74         """
75         self.train_flag = False
76
77     def relu(self,x):
78         """
79         Relu(x) = x, x > 0
80                 0, x <= 0
81         """
82         return np.maximum(0,x)
83
84     def d_relu(self,x):
85         x[x <= 0] = 0
86         x[x > 0] = 1
87         return x
88
89     def sigmoid(self,x):
90         """
91         Element-wise function
92         \Sigma(x) = 1/(1+\ee^{-x})
93         """
94         return 1 / (1 + np.exp(-x))
95
96     def d_sigmoid(self,x):
97         """
98         Derivative of sigmoid function
99         \Sigma'(x) = \Sigma(x) * (1 - \Sigma(x))
100        """
101        return self.sigmoid(x) * (1 - self.sigmoid(x))
102
103     def tanh(self,x):
104         return np.tanh(x)
105
106     def d_tanh(self,x):
107         return 1 - np.tanh(x) ** 2
108
109     def MSE(self,y_hat,y):
110         """
111         Mean-square error (MSE)

```



```

112     """
113     return np.linalg.norm(y_hat - y) # 2-norm
114
115 def cross_entropy(self,y_hat,y):
116     """
117     Cross entropy loss
118     """
119     return y * np.log(y_hat) + (1 - y) * np.log(1 - y_hat)
120
121 def forward(self,x):
122     """
123     w/o activation:  $z^{(l+1)} = W^{(l)}a^{(l)} + b^{(l)}$ 
124     w/ activation :  $a^{(l+1)} = f(z^{(l+1)})$ 
125     """
126     # training
127     if self.train_flag:
128         self.memory["a0"] = np.copy(x)
129         x = self.fc1(x) # N * hidden
130         self.memory["z1"] = np.copy(x)
131         x = self.sigmoid(x)
132         self.memory["a1"] = np.copy(x)
133         x = self.fc2(x) # N * out
134         self.memory["z2"] = np.copy(x)
135         x = self.sigmoid(x)
136     # inferencing
137     else:
138         x = self.fc1(x) # N * hidden
139         x = self.sigmoid(x)
140         x = self.fc2(x) # N * out
141         x = self.sigmoid(x)
142     return x
143
144 def backward(self,y_hat,y,lamb=0):
145     """
146     Use Mean-Squared Error (MSE) as error function
147
148     lambda is used for weight decay
149
150     Ref: http://ufldl.stanford.edu/tutorial/supervised/MultiLayerNeuralNetworks/
151     """
152     batch_size = y.shape[0]
153     # Calculate \delta
154     # output layer:  $\delta_{(n_1)} = -(y - a_{(n_1)}) * f'(z_{(n_1)})$ 
155     # other layers:  $\delta_{(l)} = W_{(l)}^T \delta_{(l+1)} * f'(z_{(l)})$ 
156     delta = [0] * 3
157     delta[2] = (y_hat - y) * self.d_sigmoid(self.memory["z2"]) # N * out_features
158     delta[1] = np.dot(delta[2],self.fc2.weight) * self.d_sigmoid(self.memory["z1"]) # N
159     # print(delta[2].shape,delta[1].shape)
160
161     # Calculat \nabla
162     # output layer:  $\nabla_{W_{(1)}} J(W,b;x,y) = \delta_{(l+1)}(a_{(l)})^T$  # outer product
163     # other layers:  $\nabla_{b_{(1)}} J(W,b;x,y) = \delta_{(l+1)}$ 
164     nabla_W = [0] * 2
165     nabla_W[1] = np.einsum("ij,ik->ijk",delta[2],self.memory["a1"]) # N * out_features
166     # print(nabla_W[1].shape)

```

```

166 nabla_W[0] = np.einsum("ij,ik->ijk",delta[1],self.memory["a0"]) # N *
    ↪ hidden_features * in_features
167 nabla_b = [0] * 2
168 nabla_b[1] = delta[2] # N * out_features
169 nabla_b[0] = delta[1] # N * hidden_features
170 # print(nabla_W[1].shape,nabla_W[0].shape,nabla_b[1].shape,nabla_b[0].shape)
171
172 # Update parameters
173 #  $W(1) = W(1) - \alpha((1/m \Delta W(1)) + \lambda W(1))$ 
174 #  $b(1) = b(1) - \alpha(1/m \Delta b(1))$ 
175 # Use einsum to accelerate
176 # https://rockt.github.io/2018/04/30/einsum
177 nabla_W[1] = nabla_W[1].mean(axis=0)
178 nabla_W[0] = nabla_W[0].mean(axis=0)
179 nabla_b[1] = nabla_b[1].mean(axis=0)
180 nabla_b[0] = nabla_b[0].mean(axis=0)
181
182 # weight decay, lambda is the L2 regularization term
183 self.fc2.weight -= self.learning_rate * (nabla_W[1] + lamb * self.fc2.weight /
    ↪ batch_size)
184 self.fc1.weight -= self.learning_rate * (nabla_W[0] + lamb * self.fc1.weight /
    ↪ batch_size)
185 self.fc2.bias -= self.learning_rate * nabla_b[1]
186 self.fc1.bias -= self.learning_rate * nabla_b[0]

```