

## optim.py

```

1   import numpy as np
2
3   """
4   This code was originally written for CS 231n at Stanford University
5   (cs231n.stanford.edu). It has been modified in various areas for use in the
6   ECE 239AS class at UCLA. This includes the descriptions of what code to
7   implement as well as some slight potential changes in variable names to be
8   consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
9   permission to use this code. To see the original version, please visit
10  cs231n.stanford.edu.
11  """
12
13  """
14  This file implements various first-order update rules that are commonly used for
15  training neural networks. Each update rule accepts current weights and the
16  gradient of the loss with respect to those weights and produces the next set of
17  weights. Each update rule has the same interface:
18
19  def update(w, dw, config=None):
20
21  Inputs:
22      - w: A numpy array giving the current weights.
23      - dw: A numpy array of the same shape as w giving the gradient of the
24          loss with respect to w.
25      - config: A dictionary containing hyperparameter values such as learning rate,
26          momentum, etc. If the update rule requires caching values over many
27          iterations, then config will also hold these cached values.
28
29  Returns:
30      - next_w: The next point after the update.
31      - config: The config dictionary to be passed to the next iteration of the
32          update rule.
33
34  NOTE: For most update rules, the default learning rate will probably not perform
35  well; however the default values of the other hyperparameters should work well
36  for a variety of different problems.
37
38  For efficiency, update rules may perform in-place updates, mutating w and
39  setting next_w equal to w.
40  """
41
42
43  def sgd(w, dw, config=None):
44  """
45      Performs vanilla stochastic gradient descent.
46
47      config format:
48          - learning_rate: Scalar learning rate.
49  """
50      if config is None: config = {}
51      config.setdefault('learning_rate', 1e-2)
52
53      w -= config['learning_rate'] * dw
54      return w, config
55
56
57  def sgd_momentum(w, dw, config=None):
58  """
59      Performs stochastic gradient descent with momentum.
60
61      config format:
62          - learning_rate: Scalar learning rate.
63          - momentum: Scalar between 0 and 1 giving the momentum value.
64              Setting momentum = 0 reduces to sgd.
65          - velocity: A numpy array of the same shape as w and dw used to store a moving
66              average of the gradients.
67  """
68      if config is None: config = {}
69      config.setdefault('learning_rate', 1e-2)
70      config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
71      v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

```

```

72     # ===== #
73     # YOUR CODE HERE:
74     #   Implement the momentum update formula.  Return the updated weights
75     #   as next_w, and the updated velocity as v.
76     # ===== #
77     v = config['momentum'] * v - config['learning_rate'] * dw
78     next_w = w + v
79     # ===== #
80     # END YOUR CODE HERE
81     # ===== #
82
83     config['velocity'] = v
84
85     return next_w, config
86
87
88 def sgd_nesterov_momentum(w, dw, config=None):
89     """
90     Performs stochastic gradient descent with Nesterov momentum.
91
92     config format:
93     - learning_rate: Scalar learning rate.
94     - momentum: Scalar between 0 and 1 giving the momentum value.
95     Setting momentum = 0 reduces to sgd.
96     - velocity: A numpy array of the same shape as w and dw used to store a moving
97     average of the gradients.
98     """
99     if config is None: config = {}
100    config.setdefault('learning_rate', 1e-2)
101    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
102    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
103
104    # ===== #
105    # YOUR CODE HERE:
106    #   Implement the momentum update formula.  Return the updated weights
107    #   as next_w, and the updated velocity as v.
108    # ===== #
109    v_old = v
110    v = config['momentum'] * v - config['learning_rate'] * dw
111    next_w = w + v + config['momentum'] * (v - v_old)
112    # ===== #
113    # END YOUR CODE HERE
114    # ===== #
115
116    config['velocity'] = v
117
118    return next_w, config
119
120
121 def rmsprop(w, dw, config=None):
122     """
123     Uses the RMSProp update rule, which uses a moving average of squared gradient
124     values to set adaptive per-parameter learning rates.
125
126     config format:
127     - learning_rate: Scalar learning rate.
128     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
129     gradient cache.
130     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
131     - beta: Moving average of second moments of gradients.
132     """
133     if config is None: config = {}
134     config.setdefault('learning_rate', 1e-2)
135     config.setdefault('decay_rate', 0.99)
136     config.setdefault('epsilon', 1e-8)
137     config.setdefault('a', np.zeros_like(w))
138
139     next_w = None
140
141    # ===== #
142    # YOUR CODE HERE:
143    #   Implement RMSProp.  Store the next value of w as next_w.  You need
144    #   to also store in config['a'] the moving average of the second
145    #   moment gradients, so they can be used for future gradients. Concretely,
146    #   config['a'] corresponds to "a" in the lecture notes.

```

```

146 # ===== #
147 config['a'] = config['decay_rate'] * config['a'] + (1-config['decay_rate']) * np.multiply(dw, dw)
148 next_w = w - np.multiply(config['learning_rate'] / (np.sqrt(config['a']) + config['epsilon']), dw)
149 # ===== #
150 # END YOUR CODE HERE
151 # ===== #
152
153 return next_w, config
154
155
156 def adam(w, dw, config=None):
157 """
158 Uses the Adam update rule, which incorporates moving averages of both the
159 gradient and its square and a bias correction term.
160
161 config format:
162 - learning_rate: Scalar learning rate.
163 - beta1: Decay rate for moving average of first moment of gradient.
164 - beta2: Decay rate for moving average of second moment of gradient.
165 - epsilon: Small scalar used for smoothing to avoid dividing by zero.
166 - m: Moving average of gradient.
167 - v: Moving average of squared gradient.
168 - t: Iteration number.
169 """
170 if config is None: config = {}
171 config.setdefault('learning_rate', 1e-3)
172 config.setdefault('beta1', 0.9)
173 config.setdefault('beta2', 0.999)
174 config.setdefault('epsilon', 1e-8)
175 config.setdefault('v', np.zeros_like(w))
176 config.setdefault('a', np.zeros_like(w))
177 config.setdefault('t', 0)
178
179 next_w = None
180
181 # ===== #
182 # YOUR CODE HERE:
183 # Implement Adam. Store the next value of w as next_w. You need
184 # to also store in config['a'] the moving average of the second
185 # moment gradients, and in config['v'] the moving average of the
186 # first moments. Finally, store in config['t'] the increasing time.
187 # ===== #
188 config['t'] += 1
189 config['v'] = config['beta1'] * config['v'] + (1-config['beta1']) * dw
190 config['a'] = config['beta2'] * config['a'] + (1-config['beta2']) * np.multiply(dw, dw)
191 v_corr = config['v']/(1-(config['beta1'] ** config['t']))
192 a_corr = config['a']/(1-(config['beta2'] ** config['t']))
193 next_w = w - np.multiply((config['learning_rate']/(np.sqrt(a_corr) + config['epsilon'])), v_corr)
194 # ===== #
195 # END YOUR CODE HERE
196 # ===== #
197
198 return next_w, config
199

```