

```
1 import numpy as np
2
3
4 """
5 This file implements various first-order update rules that are commonly used
6 for
7 training neural networks. Each update rule accepts current weights and the
8 gradient of the loss with respect to those weights and produces the next set
9 of
10 weights. Each update rule has the same interface:
11
12 def update(w, dw, config=None):
13
14 Inputs:
15 - w: A numpy array giving the current weights.
16 - dw: A numpy array of the same shape as w giving the gradient of the
17 loss with respect to w.
18 - config: A dictionary containing hyperparameter values such as learning
19 rate,
20 momentum, etc. If the update rule requires caching values over many
21 iterations, then config will also hold these cached values.
22
23 Returns:
24 - next_w: The next point after the update.
25 - config: The config dictionary to be passed to the next iteration of the
26 update rule.
27
28 NOTE: For most update rules, the default learning rate will probably not
29 perform
30 well; however the default values of the other hyperparameters should work
31 well
32 for a variety of different problems.
33
34 For efficiency, update rules may perform in-place updates, mutating w and
35 setting next_w equal to w.
36 """
37
38 def sgd(w, dw, config=None):
39     """
40     Performs vanilla stochastic gradient descent.
41
42     config format:
43     - learning_rate: Scalar learning rate.
44     """
45     if config is None: config = {}
46     config.setdefault('learning_rate', 1e-2)
47
48     w -= config['learning_rate'] * dw
49     return w, config
50
51 def sgd_momentum(w, dw, config=None):
52     """
53     Performs stochastic gradient descent with momentum.
54
55     config format:
56     - learning_rate: Scalar learning rate.
57     - momentum: Scalar between 0 and 1 giving the momentum value.
```

```

55     Setting momentum = 0 reduces to sgd.
56     - velocity: A numpy array of the same shape as w and dw used to store a
moving
57     average of the gradients.
58     """
59     if config is None: config = {}
60     config.setdefault('learning_rate', 1e-2)
61     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
62     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
63
64     # ===== #
65     # YOUR CODE HERE:
66     # Implement the momentum update formula. Return the updated weights
67     # as next_w, and the updated velocity as v.
68     # ===== #
69
70     alpha = config['momentum']
71     eps = config['learning_rate']
72
73     v = alpha * v - eps * dw
74     w += v
75
76     next_w = w
77
78     # ===== #
79     # END YOUR CODE HERE
80     # ===== #
81
82     config['velocity'] = v
83
84     return next_w, config
85
86 def sgd_nesterov_momentum(w, dw, config=None):
87     """
88     Performs stochastic gradient descent with Nesterov momentum.
89
90     config format:
91     - learning_rate: Scalar learning rate.
92     - momentum: Scalar between 0 and 1 giving the momentum value.
93     Setting momentum = 0 reduces to sgd.
94     - velocity: A numpy array of the same shape as w and dw used to store a
moving
95     average of the gradients.
96     """
97     if config is None: config = {}
98     config.setdefault('learning_rate', 1e-2)
99     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
100    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
101
102    # ===== #
103    # YOUR CODE HERE:
104    # Implement the momentum update formula. Return the updated weights
105    # as next_w, and the updated velocity as v.
106    # ===== #
107
108    alpha = config['momentum']
109    eps = config['learning_rate']
110

```

```

111     v_old = v
112     v = alpha * v - eps * dw
113     w += (v + alpha * (v - v_old))
114
115     next_w = w
116
117     # ===== #
118     # END YOUR CODE HERE
119     # ===== #
120
121     config['velocity'] = v
122
123     return next_w, config
124
125 def rmsprop(w, dw, config=None):
126     """
127     Uses the RMSProp update rule, which uses a moving average of squared
128     gradient
129     values to set adaptive per-parameter learning rates.
130
131     config format:
132     - learning_rate: Scalar learning rate.
133     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
134     gradient cache.
135     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
136     - beta: Moving average of second moments of gradients.
137     """
138     if config is None: config = {}
139     config.setdefault('learning_rate', 1e-2)
140     config.setdefault('decay_rate', 0.99)
141     config.setdefault('epsilon', 1e-8)
142     config.setdefault('a', np.zeros_like(w))
143
144     next_w = None
145
146     # ===== #
147     # YOUR CODE HERE:
148     # Implement RMSProp. Store the next value of w as next_w. You need
149     # to also store in config['a'] the moving average of the second
150     # moment gradients, so they can be used for future gradients. Concretely,
151     # config['a'] corresponds to "a" in the lecture notes.
152     # ===== #
153
154     a = config['a']
155     beta = config['decay_rate']
156     eps = config['learning_rate']
157     nu = config['epsilon']
158
159     a = beta * a + (1 - beta) * dw * dw
160     w -= (eps * dw) / (np.sqrt(a) + nu)
161
162     config['a'] = a
163     next_w = w
164
165     # ===== #
166     # END YOUR CODE HERE
167     # ===== #
168
169     return next_w, config

```

```

170
171 def adam(w, dw, config=None):
172     """
173     Uses the Adam update rule, which incorporates moving averages of both the
174     gradient and its square and a bias correction term.
175
176     config format:
177     - learning_rate: Scalar learning rate.
178     - beta1: Decay rate for moving average of first moment of gradient.
179     - beta2: Decay rate for moving average of second moment of gradient.
180     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
181     - m: Moving average of gradient.
182     - v: Moving average of squared gradient.
183     - t: Iteration number.
184     """
185     if config is None: config = {}
186     config.setdefault('learning_rate', 1e-3)
187     config.setdefault('beta1', 0.9)
188     config.setdefault('beta2', 0.999)
189     config.setdefault('epsilon', 1e-8)
190     config.setdefault('v', np.zeros_like(w))
191     config.setdefault('a', np.zeros_like(w))
192     config.setdefault('t', 0)
193
194     next_w = None
195
196     # ===== #
197     # YOUR CODE HERE:
198     # Implement Adam. Store the next value of w as next_w. You need
199     # to also store in config['a'] the moving average of the second
200     # moment gradients, and in config['v'] the moving average of the
201     # first moments. Finally, store in config['t'] the increasing time.
202     # ===== #
203
204     t = config['t']
205     v = config['v']
206     a = config['a']
207     eps = config['learning_rate']
208     nu = config['epsilon']
209     beta1 = config['beta1']
210     beta2 = config['beta2']
211
212     t += 1
213     v = beta1 * v + (1 - beta1) * dw
214     a = beta2 * a + (1 - beta2) * dw * dw
215     v_u = v / (1 - beta1**t)
216     a_u = a / (1 - beta2**t)
217     w -= (eps * v_u) / (np.sqrt(a_u) + nu)
218
219     config['t'] = t
220     config['v'] = v
221     config['a'] = a
222     next_w = w
223
224     # ===== #
225     # END YOUR CODE HERE
226     # ===== #
227
228     return next_w, config
229

```

230
231
232
233
234