```python
import numpy as np

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu).  It has been modified in various areas for use in the
ECE 239AS class at UCLA.  This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung
for
permission to use this code.  To see the original version, please visit
cs231n.stanford.edu.
"""

"""
This file implements various first-order update rules that are commonly used
for
training neural networks. Each update rule accepts current weights and the
gradient of the loss with respect to those weights and produces the next set
of
weights. Each update rule has the same interface:

def update(w, dw, config=None):

Inputs:
  - w: A numpy array giving the current weights.
  - dw: A numpy array of the same shape as w giving the gradient of the
    loss with respect to w.
  - config: A dictionary containing hyperparameter values such as learning
rate,
    momentum, etc. If the update rule requires caching values over many
    iterations, then config will also hold these cached values.

Returns:
  - next_w: The next point after the update.
  - config: The config dictionary to be passed to the next iteration of the
    update rule.

NOTE: For most update rules, the default learning rate will probably not
perform
well; however the default values of the other hyperparameters should work
well
for a variety of different problems.

For efficiency, update rules may perform in-place updates, mutating w and
setting next_w equal to w.
"""


def sgd(w, dw, config=None):
    """
    Performs vanilla stochastic gradient descent.

    config format:
    - learning_rate: Scalar learning rate.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)

    w -= config['learning_rate'] * dw
    return w, config
```

```python
def sgd_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets it to zero.

    # ================================================================ #
    # YOUR CODE HERE:
    #   Implement the momentum update formula.  Return the updated weights
    #   as next_w, and the updated velocity as v.
    # ================================================================ #

    alpha = config.get('momentum')
    eps = config.get('learning_rate')
    v = alpha*v - eps*dw
    next_w = w+v

    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    config['velocity'] = v

    return next_w, config

def sgd_nesterov_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with Nesterov momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets it to zero.

    # ================================================================ #
```

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

```python
    # =============================================================== #
    # END YOUR CODE HERE
    # =============================================================== #

    return next_w, config


def adam(w, dw, config=None):
    """
    Uses the Adam update rule, which incorporates moving averages of both the
    gradient and its square and a bias correction term.

    config format:
    - learning_rate: Scalar learning rate.
    - beta1: Decay rate for moving average of first moment of gradient.
    - beta2: Decay rate for moving average of second moment of gradient.
    - epsilon: Small scalar used for smoothing to avoid dividing by zero.
    - m: Moving average of gradient.
    - v: Moving average of squared gradient.
    - t: Iteration number.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-3)
    config.setdefault('beta1', 0.9)
    config.setdefault('beta2', 0.999)
    config.setdefault('epsilon', 1e-8)
    config.setdefault('v', np.zeros_like(w))
    config.setdefault('a', np.zeros_like(w))
    config.setdefault('t', 0)

    next_w = None

    # =============================================================== #
    # YOUR CODE HERE:
    #   Implement Adam.  Store the next value of w as next_w.  You need
    #   to also store in config['a'] the moving average of the second
    #   moment gradients, and in config['v'] the moving average of the
    #   first moments.  Finally, store in config['t'] the increasing time.
    # =============================================================== #

    t = config.get('t') + 1
    v = config.get('v')
    beta1 = config.get('beta1')
    v = beta1*v + (1-beta1)*dw

    a = config.get('a')
    beta2 = config.get('beta2')
    a = beta2*a + (1-beta2)*(dw**2)

    v_u = v / (1 - beta1**t)
    a_u = a / (1 - beta2**t)

    lr = config.get('learning_rate')
    eps = config.get('epsilon')
    next_w = w - (lr / (np.sqrt(a_u) + eps) * v_u)

    config['a'] = a
    config['v'] = v
    config['t'] = t
```

```python
    # ================================================================= #
    # END YOUR CODE HERE
    # ================================================================= #

    return next_w, config
```