

Note that the iterator i starts at 1 in the `for` loop as the first parameters are $W^{(1)}$ and $b^{(1)}$.

$$\begin{aligned} \hat{r}_{\text{Adam}}^{(i)} &= \hat{r}_{\text{Adam}}^{(i-1)} + (1 - \beta_1) \frac{g^{(i)}}{\sqrt{m^{(i)}}} \\ \hat{r}_{\text{Adam}}^{(i)} &= \frac{\hat{r}_{\text{Adam}}^{(i-1)}}{1 - \beta_1^2} \\ \hat{r}_{\text{Adam}}^{(i)} &= \hat{r}_{\text{Adam}}^{(i-1)} + (1 - \beta_1) \frac{g^{(i)}}{\sqrt{m^{(i)}}} \\ \hat{r}_{\text{Adam}}^{(i)} &= \frac{\hat{r}_{\text{Adam}}^{(i-1)}}{1 - \beta_1^2} \\ \hat{r}_{\text{Adam}}^{(i)} &= \hat{r}_{\text{Adam}}^{(i-1)} + (1 - \beta_1) \frac{g^{(i)}}{\sqrt{m^{(i)}}} \\ \hat{r}_{\text{Adam}}^{(i)} &= \frac{\hat{r}_{\text{Adam}}^{(i-1)}}{1 - \beta_1^2} \end{aligned}$$

```
In [14]: # GRADED FUNCTION: update_parameters_with_adam

def update_parameters_with_adam(parameters, grads, v, s, learning_rate = 0.01,
                                beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8):
    """Update parameters using Adam

    Arguments:
    parameters -- python dictionary containing your parameters:
        parameters["W"] = w1
        parameters["b"] = b1
    grads -- python dictionary containing your gradients for each parameters:
        grads["dw"] = dw1
        grads["db"] = db1
    v -- Adam variable, moving average of the first gradient, python dictionary
    s -- Adam variable, moving average of the squared gradient, python dictionary
    learning_rate -- the learning rate, scalar.
    beta1 -- Exponential decay hyperparameter for the first moment estimates
    beta2 -- Exponential decay hyperparameter for the second moment estimates
    epsilon -- hyperparameter preventing division by zero in Adam updates

    Returns:
    parameters -- python dictionary containing your updated parameters
    """
    # Initialize first moment estimate, python dictionary
    v = {}
    # Initialize second moment estimate, python dictionary
    s = {}

    # Iterate Adam update on all parameters
    for i in range(1, L + 1):
        # Moving average of the gradients. Inputs: "v, grads, beta1". Output: "v".
        # (approx. 2 lines)
        v["dw"] = beta1 * v["dw"] + (1 - beta1) * grads["dw"]
        v["db"] = beta1 * v["db"] + (1 - beta1) * grads["db"]
        # YOUR CODE STARTS HERE
        # Update parameters with Adam
        parameters["W"] = parameters["W"] - learning_rate * v["dw"] / (1 - beta1 + epsilon)
        parameters["b"] = parameters["b"] - learning_rate * v["db"] / (1 - beta1 + epsilon)
        # YOUR CODE ENDS HERE

    # Compute bias-corrected first moment estimate. Inputs: "v, beta1, t". Output: "v_corrected".
    # (approx. 2 lines)
    v_corrected["dw"] = v["dw"] / (1 - beta1 ** t)
    v_corrected["db"] = v["db"] / (1 - beta1 ** t)
    # YOUR CODE STARTS HERE
    # Update parameters with Adam
    parameters["W"] = parameters["W"] - learning_rate * v_corrected["dw"] / (1 - beta1 + epsilon)
    parameters["b"] = parameters["b"] - learning_rate * v_corrected["db"] / (1 - beta1 + epsilon)
    # YOUR CODE ENDS HERE

    return parameters, v, s, v_corrected, s_corrected, epsilon

# Update parameters. Inputs: "parameters, learning_rate, v_corrected, s_corrected, epsilon". Output: "parameters"
# (approx. 2 lines)
parameters["W"] = parameters["W"] - learning_rate * v_corrected["dw"] / (1 - beta1 + epsilon)
parameters["b"] = parameters["b"] - learning_rate * v_corrected["db"] / (1 - beta1 + epsilon)
# YOUR CODE ENDS HERE

return parameters, v, s, v_corrected, s_corrected, epsilon
```

```
In [15]: parameters, grads, v, s = update_parameters_with_adam_test_case()

t = 2
learning_rate = 0.02
beta1 = 0.9
beta2 = 0.999
epsilon = 1e-8

parameters, v, s, v_corrected, s_corrected, epsilon = update_parameters_with_adam_test_case(parameters, grads, v, s, t, learning_rate, beta1, beta2, epsilon)

print("W1 = {}".format(parameters["W1"]))
print("W2 = {}".format(parameters["W2"]))
print("b1 = {}".format(parameters["b1"]))
print("b2 = {}".format(parameters["b2"]))

update_parameters_with_adam_test(update_parameters_with_adam)

W1 =
[[ 1.6392428 -0.6268425 -0.54320974]
 [ 1.0870243  0.8503893 -0.2865723]]
W2 =
[[ 0.2556139 -0.2442599  1.4770772]
 [-0.4538458 -0.2674493 -0.3469034]]
b1 =
[[ 1.14873036 -1.09256871 -0.15734651]]
b2 =
[[ 1.75854357]
 [-0.7456607]]
epsilon =
[[ 0.89228024]
 [ 0.89270193]]
All test passed

Expected values:
W1 =
[[ 1.6392428 -0.6268425 -0.54320974]
 [ 1.0870243  0.8503893 -0.2865723]]
W2 =
[[ 0.2556139 -0.2442599  1.4770772]
 [-0.4538458 -0.2674493 -0.3469034]]
b1 =
[[ 1.14873036 -1.09256871 -0.15734651]]
b2 =
[[ 1.75854357]
 [-0.7456607]]
epsilon =
[[ 0.89228024]
 [ 0.89270193]]
All test passed
```

You now have three working optimization algorithms (mini-batch gradient descent, Momentum, Adam). Let's implement a model with each of these optimizers and observe the difference.

6 - Model with different Optimization algorithms

Below, you'll use a `mnist` dataset to test the different optimization methods. (The dataset is named "mnist" because the data from each of the two classes looks a bit like a crescent-shaped moon.)

```
In [27]: train_X, train_Y = load_dataset()

# Plot decision boundary
plt.figure(figsize=(10, 10))
plt.scatter(train_X, train_Y)
plt.title('Decision boundary')
plt.show()
```

A 3-layer neural network has already been implemented for you. You'll test it with:

- Mini-batch Gradient Descent: it will call your function:
- Mini-batch Momentum: it will call your function:
- Mini-batch Adam: it will call your function:

```
In [28]: def model(X, Y, layers_dims, optimizer, learning_rate = 0.0007, mini_batch_size = 64, beta = 0.9,
                beta2 = 0.999, epsilon = 1e-8, num_epochs = 5000, print_cost = True):
    """3-layer neural network model which can be run in different optimizer modes.

    Arguments:
    X -- input data, of shape (2, number of examples)
    Y -- true "label" vector (1 for blue dot / 0 for red dot), of shape (1, number of examples)
    layers_dims -- python list, containing the size of each layer
    learning_rate -- the learning rate, scalar.
    mini_batch_size -- the size of a mini batch
    beta -- Momentum hyperparameter
    beta2 -- Exponential decay hyperparameter for the past gradients estimates
    epsilon -- Exponential decay hyperparameter for the past squared gradients estimates
    num_epochs -- number of epochs
    print_cost -- True to print the cost every 1000 epochs

    Returns:
    parameters -- python dictionary containing your updated parameters
    """
    L = len(layers_dims) # number of layers in the neural networks
    costs = [] # to keep track of the cost
    t = 0 # Initializing the counter required for Adam update
    seed = 10 # For grading purposes, so that your "random" minibatches are the same as ours
    m = X.shape[1] # number of training examples

    # Initialize parameters
    parameters = initialize_parameters(layers_dims)

    # Initialize the optimizer
    if optimizer == "gd":
        pass # no initialization required for gradient descent
    elif optimizer == "momentum":
        v = initialize_velocity(parameters)
    elif optimizer == "adam":
        v, s = initialize_adam(parameters)

    # Optimization loop
    for i in range(num_epochs):
        # Define the random minibatches. We increment the seed to reshuffle differently the dataset after each epoch
        seed = seed + 1
        minibatches = random_mini_batches(X, Y, mini_batch_size, seed)
        cost_total = 0

        for minibatch in minibatches:
            # Select a minibatch
            (minibatch_X, minibatch_Y) = minibatch

            # Forward propagation
            a1, caches = forward_propagation(minibatch_X, parameters)

            # Compute cost and add to the cost total
            cost_total += compute_cost(a1, minibatch_Y)

            # Backward propagation
            grads = backward_propagation(minibatch_X, minibatch_Y, caches)

            # Update parameters
            if optimizer == "gd":
                parameters = update_parameters_with_gd(parameters, grads, learning_rate)
            elif optimizer == "momentum":
                parameters, v = update_parameters_with_momentum(parameters, grads, v, beta, learning_rate)
            elif optimizer == "adam":
                parameters, v, s, _ = update_parameters_with_adam(parameters, grads, v, s,
                                                                    beta, learning_rate, beta2, epsilon)

        cost_avg = cost_total / m

        # Print the cost every 1000 epoch
        if print_cost and i % 1000 == 0:
            print ("Cost after epoch %i: %f" % i, cost_avg)
        if print_cost and i % 100 == 0:
            costs.append(cost_avg)

    # plot the costs
    plt.plot(costs)
    plt.ylabel('cost')
    plt.xlabel('epoch (per 100)')
    plt.title('Learning rate = ' + str(learning_rate))
    plt.show()

    return parameters
```

Now, run this 3 layer neural network with each of the 3 optimization methods.

6.1 - Mini-Batch Gradient Descent

Run the following code to see how the model does with mini-batch gradient descent.

```
In [29]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "gd")

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title('Model with Gradient Descent optimization')
axes = plt.gca()
axes.set_xlim([-1.5, 2.5])
axes.set_ylim([-1, 1])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)

Cost after epoch 0: 0.702405
Cost after epoch 1000: 0.668101
Cost after epoch 2000: 0.653588
Cost after epoch 3000: 0.639788
Cost after epoch 4000: 0.630606

Learning rate = 0.0007

Accuracy: 0.7166666666666667

Model with Gradient Descent optimization
```

6.2 - Mini-Batch Gradient Descent with Momentum

Next, run the following code to see how the model does with momentum. Because this example is relatively simple, the gains from using momentum are small – but for more complex problems you might see bigger gains.

```
In [30]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, beta = 0.9, optimizer = "momentum")

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title('Model with Momentum optimization')
axes = plt.gca()
axes.set_xlim([-1.5, 2.5])
axes.set_ylim([-1, 1])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)

Cost after epoch 0: 0.702405
Cost after epoch 1000: 0.661881
Cost after epoch 2000: 0.653588
Cost after epoch 3000: 0.639788
Cost after epoch 4000: 0.630606

Learning rate = 0.0007

Accuracy: 0.7166666666666667

Model with Momentum optimization
```

6.3 - Mini-Batch with Adam

Finally, run the following code to see how the model does with Adam.

```
In [31]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "adam")

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title('Model with Adam optimization')
axes = plt.gca()
axes.set_xlim([-1.5, 2.5])
axes.set_ylim([-1, 1])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)

Cost after epoch 0: 0.702166
Cost after epoch 1000: 0.617845
Cost after epoch 2000: 0.613116
Cost after epoch 3000: 0.638788
Cost after epoch 4000: 0.630606

Learning rate = 0.0007

Accuracy: 0.9433333333333334

Model with Adam optimization
```

6.4 - Summary

optimization method	accuracy	cost shape
Gradient descent	~71%	smooth
Momentum	~71%	smooth
Adam	~86%	smoother

Momentum usually helps, but given the small learning rate and the simple dataset, its impact is almost negligible. On the other hand, Adam clearly outperforms mini-batch gradient descent and Momentum. If you run the model for more epochs on this simple dataset, all three methods will lead to very good results. However, you've seen that Adam converges a lot faster.

Some advantages of Adam include:

- Relatively low memory requirements (though higher than gradient descent and gradient descent with momentum)
- Usually works well even with little tuning of hyperparameters (except ϵ)

References

- Adam paper: <https://arxiv.org/pdf/1412.0093v1.pdf>

7 - Learning Rate Decay and Scheduling

During the first part of training, your model can get stuck by using large step learning, but over time, using a fixed value for the learning rate alone can cause your model to get stuck in a wide oscillation that never quite converges. But if you were to slowly reduce your learning rate alpha over time, you could then take smaller, slower steps that bring you closer to the minimum. This is the idea behind learning rate decay.

Learning rate decay can be achieved by using either **adaptive methods** or **pre-defined learning rate schedules**.

Now, you'll apply scheduled learning rate decay to a 3-layer neural network in three different optimizer modes and see how each one differs, as well as the effect of scheduling at different epochs.

This model is essentially the same as the one you used before, except in this one you'll be able to include learning rate decay. It includes two new parameters, `decay` and `decay_rate`.

```
In [32]: def model(X, Y, layers_dims, optimizer, learning_rate = 0.0007, mini_batch_size = 64, beta = 0.9,
                beta2 = 0.999, epsilon = 1e-8, num_epochs = 5000, print_cost = True, decay=0.95, decay_rate=1):
    """3-layer neural network model which can be run in different optimizer modes.

    Arguments:
    X -- input data, of shape (2, number of examples)
    Y -- true "label" vector (1 for blue dot / 0 for red dot), of shape (1, number of examples)
    layers_dims -- python list, containing the size of each layer
    learning_rate -- the learning rate, scalar.
    mini_batch_size -- the size of a mini batch
    beta -- Momentum hyperparameter
    beta2 -- Exponential decay hyperparameter for the past gradients estimates
    epsilon -- Exponential decay hyperparameter for the past squared gradients estimates
    num_epochs -- number of epochs
    print_cost -- True to print the cost every 1000 epochs

    Returns:
    parameters -- python dictionary containing your updated parameters
    """
    L = len(layers_dims) # number of layers in the neural networks
    costs = [] # to keep track of the cost
    t = 0 # Initializing the counter required for Adam update
    seed = 10 # For grading purposes, so that your "random" minibatches are the same as ours
    m = X.shape[1] # number of training examples
    learning_rate = learning_rate # the original learning rate

    # Initialize parameters
    parameters = initialize_parameters(layers_dims)

    # Initialize the optimizer
    if optimizer == "gd":
        pass # no initialization required for gradient descent
    elif optimizer == "momentum":
        v = initialize_velocity(parameters)
    elif optimizer == "adam":
        v, s = initialize_adam(parameters)

    # Optimization loop
    for i in range(num_epochs):
        # Define the random minibatches. We increment the seed to reshuffle differently the dataset after each epoch
        seed = seed + 1
        minibatches = random_mini_batches(X, Y, mini_batch_size, seed)
        cost_total = 0

        for minibatch in minibatches:
            # Select a minibatch
            (minibatch_X, minibatch_Y) = minibatch

            # Forward propagation
            a1, caches = forward_propagation(minibatch_X, parameters)

            # Compute cost and add to the cost total
            cost_total += compute_cost(a1, minibatch_Y)

            # Backward propagation
            grads = backward_propagation(minibatch_X, minibatch_Y, caches)

            # Update parameters
            if optimizer == "gd":
                parameters = update_parameters_with_gd(parameters, grads, learning_rate)
            elif optimizer == "momentum":
                parameters, v = update_parameters_with_momentum(parameters, grads, v, beta, learning_rate)
            elif optimizer == "adam":
                parameters, v, s, _ = update_parameters_with_adam(parameters, grads, v, s,
                                                                    beta, learning_rate, beta2, epsilon)

        cost_avg = cost_total / m

        # Learning rate decay
        if decay:
            learning_rate = decay_rate * learning_rate

        # Print the cost every 1000 epoch
        if print_cost and i % 1000 == 0:
            print ("Cost after epoch %i: %f" % i, cost_avg)
        if print_cost and i % 100 == 0:
            costs.append(cost_avg)

    # plot the costs
    plt.plot(costs)
    plt.ylabel('cost')
    plt.xlabel('epoch (per 100)')
    plt.title('Learning rate = ' + str(learning_rate))
    plt.show()

    return parameters
```

7.1 - Decay on every iteration

For the prior of the assignment, you'll try one of the pre-defined schedules for learning rate decay, called **exponential learning rate decay**. It takes this mathematical form:

$$\alpha = \frac{1}{1 + \text{decay_Rate} \times \text{epoch_Number}^d}$$

Exercise 7 - update_lr

Calculate the new learning rate using exponential weight decay.

```
In [37]: # GRADED FUNCTION: update_lr

def update_lr(learning_rate, epoch_num, decay_rate):
    """Calculate updated learning rate using exponential weight decay.

    Arguments:
    learning_rate -- Original learning rate. Scalar
    epoch_num -- Epoch number. Integer
    decay_rate -- Decay rate. Scalar

    Returns:
    learning_rate -- Updated learning rate. Scalar
    """
    # (approx. 1 line)
    # learning_rate = ...
    # YOUR CODE STARTS HERE
    learning_rate = learning_rate * decay_rate ** (epoch_num / 1000)
    # YOUR CODE ENDS HERE
    return learning_rate

In [38]: learning_rate = 0.5
print("Original learning rate: ", learning_rate)
epoch_num = 10
decay_rate = 0.9
learning_rate_2 = update_lr(learning_rate, epoch_num, decay_rate)
print("Updated learning rate: ", learning_rate_2)

update_lr(learning_rate, epoch_num, decay_rate)

Original learning rate: 0.5
Updated learning rate after 10 epochs: 0.5
All test passed
```

```
In [39]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "gd", learning_rate = 0.1, num_epochs=5000, decay=update_lr)

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title('Model with Gradient Descent optimization')
axes = plt.gca()
axes.set_xlim([-1.5, 2.5])
axes.set_ylim([-1, 1])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)

Cost after epoch 0: 0.701931
Learning rate after epoch 0: 0.100000
Cost after epoch 1000: 0.611881
Learning rate after epoch 1000: 0.000100
Cost after epoch 2000: 0.608620
Learning rate after epoch 2000: 0.000050
Cost after epoch 3000: 0.656765
Learning rate after epoch 3000: 0.000033
Cost after epoch 4000: 0.654846
Learning rate after epoch 4000: 0.000025

Learning rate = 2e-05

Accuracy: 0.6533333333333333

Model with Gradient Descent optimization
```

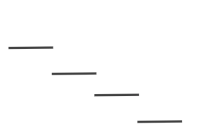
Notice that if you set the decay to occur at every iteration, the learning rate goes to zero too quickly – even if you start with a higher learning rate.

Epoch Number	Learning Rate	Cost
0	0.100000	0.701931
1000	0.000100	0.611881
2000	0.000050	0.608620
3000	0.000033	0.656765
4000	0.000025	0.654846
5000	0.000016	0.664554

When you're training for 5000 epochs it doesn't cause a lot of trouble, but when the number of epochs is large the optimization algorithm will stop updating. One common fix to this issue is to decay the learning rate every few steps. This is called **fixed interval scheduling**.

7.2 - Fixed Interval Scheduling

You can prevent the learning rate scheduling to zero too quickly by scheduling the exponential learning rate decay at a fixed time interval, for example 1000. You can either number the intervals, or divide the epoch by the time interval, which is the size of the window with the constant learning rate.



Exercise 8 - schedule_lr_decay

Calculate the new learning rate using **exponential weight decay with fixed interval scheduling**.

Instructions: Implement the learning rate scheduling such that it only changes when the `epochNum` is a multiple of the `timeInterval`.

Note: The fraction in the denominator uses the `floor` operation.

$$\alpha = \frac{1}{1 + \text{decay_Rate} \times \left\lfloor \frac{\text{epochNum}}{\text{timeInterval}} \right\rfloor^d}$$

```
In [40]: # GRADED FUNCTION: schedule_lr_decay

def schedule_lr_decay(learning_rate, epoch_num, decay_rate, time_interval=100):
    """Calculate updated learning rate using exponential weight decay.

    Arguments:
    learning_rate -- Original learning rate. Scalar
    epoch_num -- Epoch number. Integer
    decay_rate -- Decay rate. Scalar
    time_interval -- Number of epochs where you update the learning rate.

    Returns:
    learning_rate -- Updated learning rate. Scalar
    """
    # (approx. 1 line)
    # learning_rate = ...
    # YOUR CODE STARTS HERE
    learning_rate = learning_rate * decay_rate ** (epoch_num // time_interval)
    # YOUR CODE ENDS HERE
    return learning_rate

In [41]: learning_rate = 0.5
print("Original learning rate: ", learning_rate)
epoch_num = 10
decay_rate = 0.9
time_interval = 100
learning_rate_2 = schedule_lr_decay(learning_rate, epoch_num, decay_rate, time_interval)
print("Updated learning rate after 10 epochs: ", learning_rate_2)
print("Updated learning rate after 100 epochs: ", learning_rate_2)
print("Updated learning rate after 1000 epochs: ", learning_rate_2)

schedule_lr_decay(learning_rate, epoch_num, decay_rate, time_interval)

Original learning rate: 0.5
Updated learning rate after 10 epochs: 0.5
Updated learning rate after 100 epochs: 0.5
Updated learning rate after 1000 epochs: 0.38461538461538464
All test passed

Expected output:
Original learning rate: 0.5
Updated learning rate after 10 epochs: 0.5
Updated learning rate after 100 epochs: 0.5
Updated learning rate after 1000 epochs: 0.38461538461538464
```

7.3 - Using Learning Rate Decay for each Optimization Method

Below, you'll use the following "moons" dataset to test the different optimization methods. (The dataset is named "moons" because the data from each of the two classes looks a bit like a crescent-shaped moon.)

7.3.1 - Gradient Descent with Learning Rate Decay

Run the following code to see how the model does with gradient descent and weight decay.

```
In [42]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "gd", learning_rate = 0.1, num_epochs=5000, decay=update_lr)

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title('Model with Gradient Descent optimization')
axes = plt.gca()
axes.set_xlim([-1.5, 2.5])
axes.set_ylim([-1, 1])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)

Cost after epoch 0: 0.701931
Learning rate after epoch 0: 0.100000
Cost after epoch 1000: 0.611881
Learning rate after epoch 1000: 0.000100
Cost after epoch 2000: 0.608620
Learning rate after epoch 2000: 0.000050
Cost after epoch 3000: 0.656765
Learning rate after epoch 3000: 0.000033
Cost after epoch 4000: 0.654846
Learning rate after epoch 4000: 0.000025

Learning rate = 2e-05

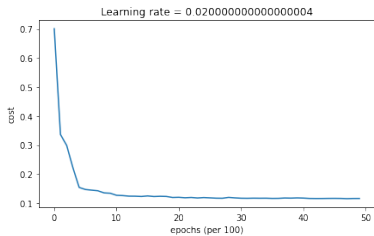
Accuracy: 0.6533333333333333

Model with Gradient Descent optimization
```

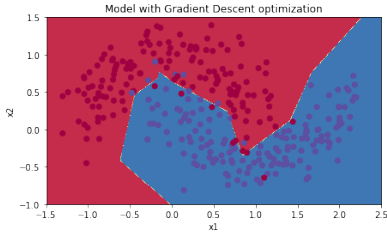
```
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title("Model with Gradient Descent optimization")
axes = plt.gca()
axes.set_xlim([-1.5,2.5])
axes.set_ylim([-1,1.5])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Cost after epoch 0: 0.701091
learning rate after epoch 0: 0.100000
Cost after epoch 1000: 0.127161
learning rate after epoch 1000: 0.050000
Cost after epoch 2000: 0.120304
learning rate after epoch 2000: 0.033333
Cost after epoch 3000: 0.117033
learning rate after epoch 3000: 0.025000
Cost after epoch 4000: 0.117512
learning rate after epoch 4000: 0.020000



Accuracy: 0.9433333333333334



7.3.2 - Gradient Descent with Momentum and Learning Rate Decay

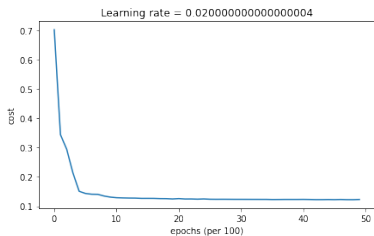
Run the following code to see how the model does gradient descent with momentum and weight decay.

```
In [43]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "momentum", learning_rate = 0.1, num_epochs=5000, decay=
schedule_lr_decay)

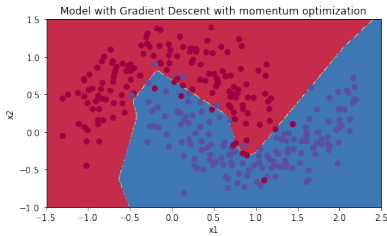
# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title("Model with Gradient Descent with momentum optimization")
axes = plt.gca()
axes.set_xlim([-1.5,2.5])
axes.set_ylim([-1,1.5])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Cost after epoch 0: 0.702226
learning rate after epoch 0: 0.100000
Cost after epoch 1000: 0.128974
learning rate after epoch 1000: 0.050000
Cost after epoch 2000: 0.125965
learning rate after epoch 2000: 0.033333
Cost after epoch 3000: 0.123375
learning rate after epoch 3000: 0.025000
Cost after epoch 4000: 0.123218
learning rate after epoch 4000: 0.020000



Accuracy: 0.9533333333333334



7.3.3 - Adam with Learning Rate Decay

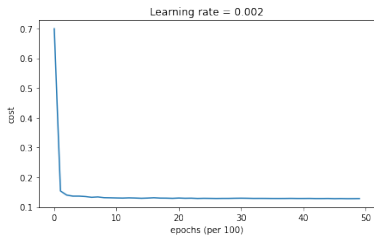
Run the following code to see how the model does Adam and weight decay.

```
In [44]: # train 3-layer model
layers_dims = [train_X.shape[0], 5, 2, 1]
parameters = model(train_X, train_Y, layers_dims, optimizer = "adam", learning_rate = 0.01, num_epochs=5000, decay=sch
edule_lr_decay)

# Predict
predictions = predict(train_X, train_Y, parameters)

# Plot decision boundary
plt.title("Model with Adam optimization")
axes = plt.gca()
axes.set_xlim([-1.5,2.5])
axes.set_ylim([-1,1.5])
plot_decision_boundary(lambda x: predict_dec(parameters, x.T), train_X, train_Y)
```

Cost after epoch 0: 0.699346
learning rate after epoch 0: 0.010000
Cost after epoch 1000: 0.130074
learning rate after epoch 1000: 0.005000
Cost after epoch 2000: 0.129826
learning rate after epoch 2000: 0.003333
Cost after epoch 3000: 0.129282
learning rate after epoch 3000: 0.002500
Cost after epoch 4000: 0.128361
learning rate after epoch 4000: 0.002000



Accuracy: 0.94



7.4 - Achieving similar performance with different methods

With Mini-batch GD or Mini-batch GD with Momentum, the accuracy is significantly lower than Adam, but when learning rate decay is added on top, either can achieve performance at a speed and accuracy score that's similar to Adam.

In the case of Adam, notice that the learning curve achieves a similar accuracy but faster.

optimization method	accuracy
Gradient descent	>94.6%
Momentum	>95.6%
Adam	94%

Congratulations! You've made it to the end of the Optimization methods notebook. Here's a quick recap of everything you're now able to do:

- Apply three different optimization methods to your models
- Build mini-batches for your training set
- Use learning rate decay scheduling to speed up your training

Great work!