

بسمه تعالی

تمرین ۲ شبکه های عصبی

سید محمدرضا حسینی

۹۷۲۴۳۱۲۹

(۱) در این تمرین یک شبکه عصبی یک لایه پنهان از پایه با استفاده از numpy نوشته شده.

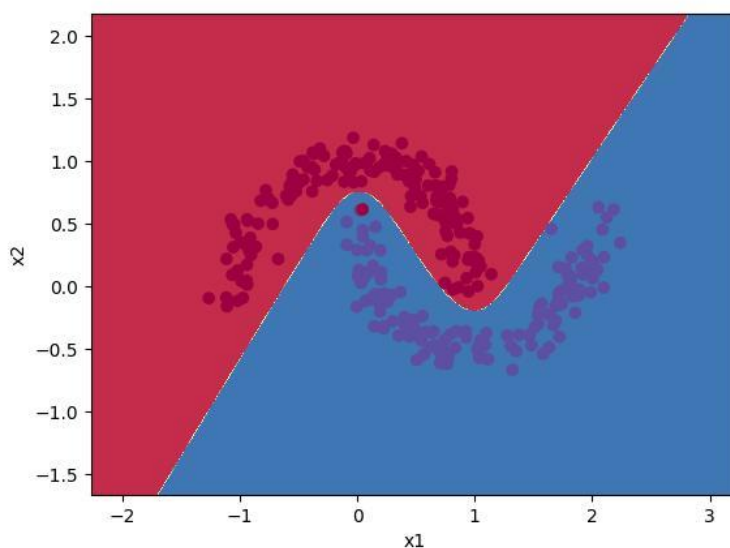
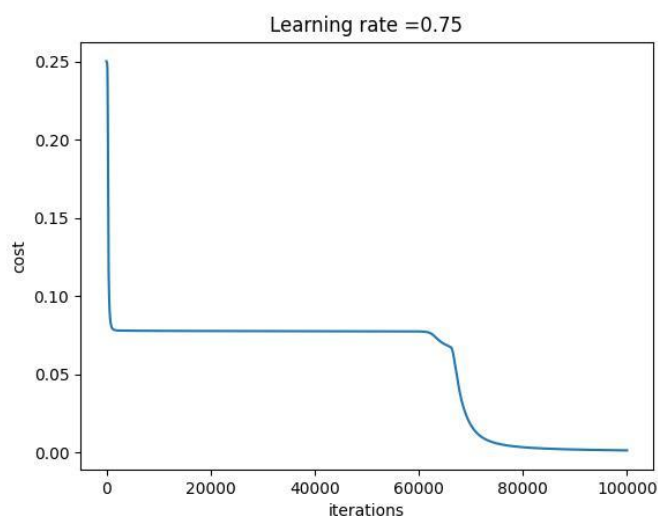
```
x, y = datasets.make_moons(n_samples=1000, shuffle=True, noise=0.1, random_state=None)
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
hidden_layer_activations = {0: "sigmoid", 1: "tanh"}
output_layer_activations = {0: "sigmoid", 1: "tanh", 2: "linear"}
hidden_activation = hidden_layer_activations.get(1) # choosing activation function of hidden layer
output_activation = output_layer_activations.get(1) # choosing activation function of hidden layer
learning_processes = {0: "Batch", 1: "Online"}
learning = learning_processes.get(0) # choosing learning process
Use_momentum = False
Use_decay_rate = True
```

در ابتدای کد متغیرهای قرار داده شده که به توان به راحتی تابع فعالیت هر لایه و نحوه آموزش را تنظیم کرد. دو متغیر بولین نیز قرار داده شده که ثابت ماندن یا تغییر کردن نرخ یادگیری و استفاده از مومنتوم در آن ست میشود.

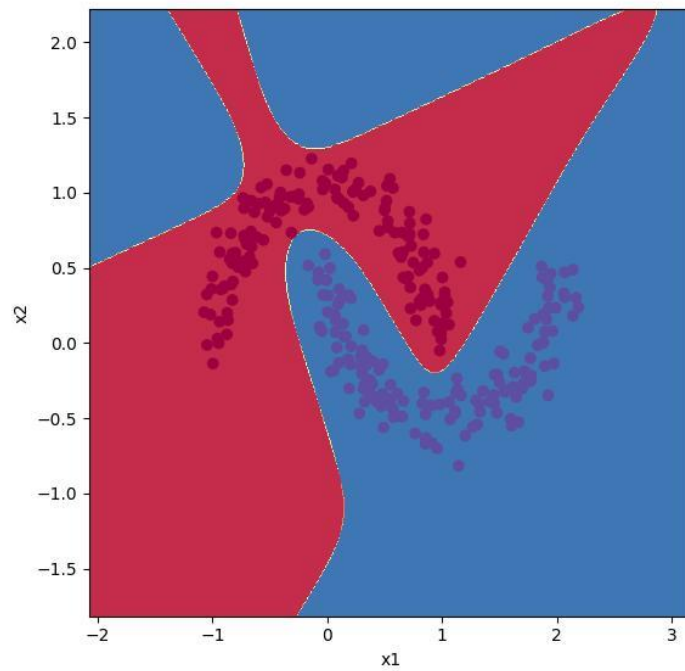
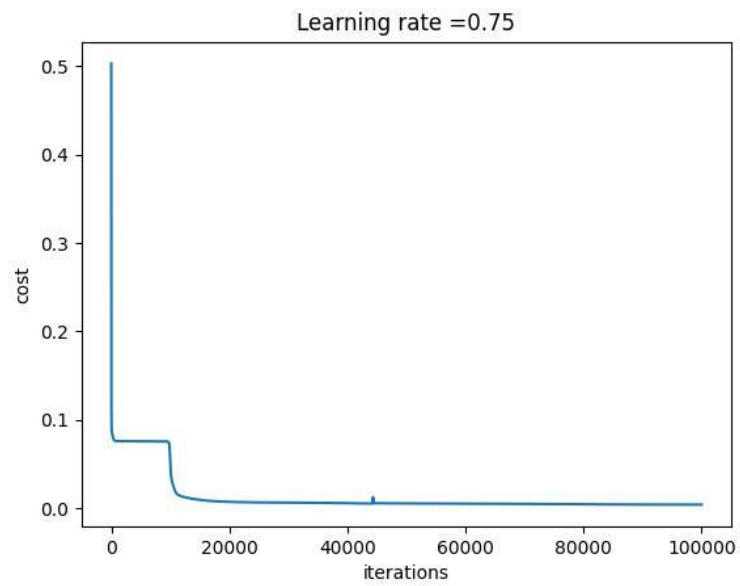
```
parameters = initialize(n_x, n_h, n_y)
for i in range(0, num_iterations):
    A2, cache = forward_propagation(X, parameters)
    cost = compute_cost(A2, Y)
    costs.append(cost)
    grads = backward_propagation(parameters, cache, X, Y)
    parameters, momentum = update_parameters(parameters, grads, learning_rate, momentum)
    if i != 0 and i % 1000 == 0 and Use_decay_rate:
        learning_rate = learning_rate * lr_decay
        # learning_rate = initial_learning_rate * (1 / (1 + lr_decay * i))
        # learning_rate = (lr_decay ** i) * initial_learning_rate
    # if cost < learning_bias:
    #     break
    if print_cost and i % 100 == 0:
        print("Cost after iteration %i: %f" % (i, cost))
    # if i > 1000:
    #     diff = 0
    #     for j in (reversed(range(len(costs) - 15, len(costs)))):
    #         diff += abs(costs[j] - costs[j - 1])
    #     if diff < 0.0001 and costs[len(costs) - 1] < 0.01:
    #         break
```

شیوه کار بدین صورت است که ابتدا مقادیر وزن ها و بایاس مقدار دهی اولیه شده و در هر حلقه بسته به نوع آموزش ابتدا forward propagation و سپس cost محاسبه میشود. بعد از آن نیز در تابع backward_propagation گرادیان نسبت به هر یک از w, b ها محاسبه میشود. در انتها نیز در تابع update_parameters نرخ w, b ها با توجه به backpropagation آپدیت میشود.

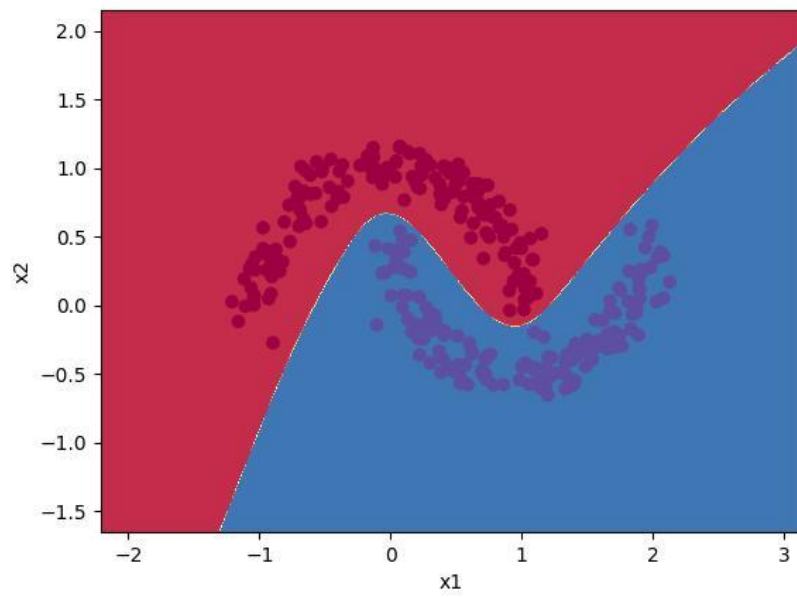
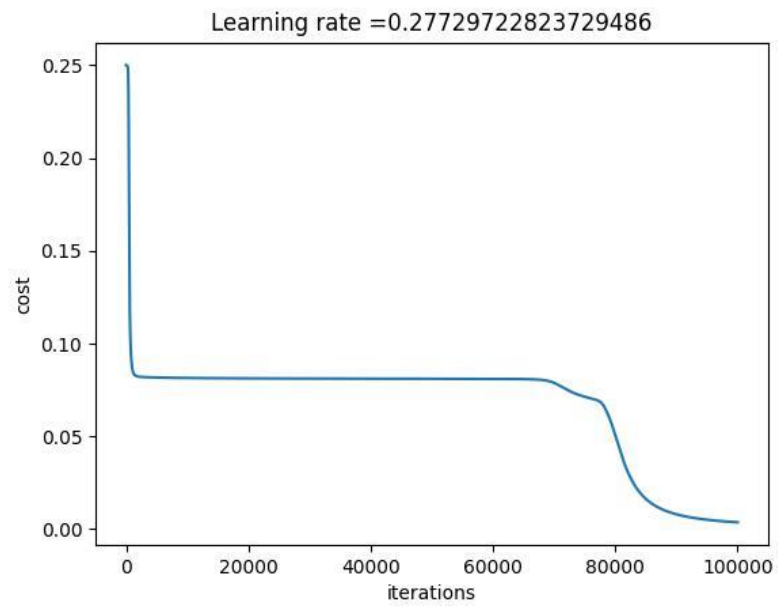
در توابع بخشی پیاده سازی شده که در صورت همگرا شدن آموزش و تغییر کم در ۱۵ epoch آموزش به اتمام برسد ولی با توجه به توضیحات این بخش کامنت شده است.



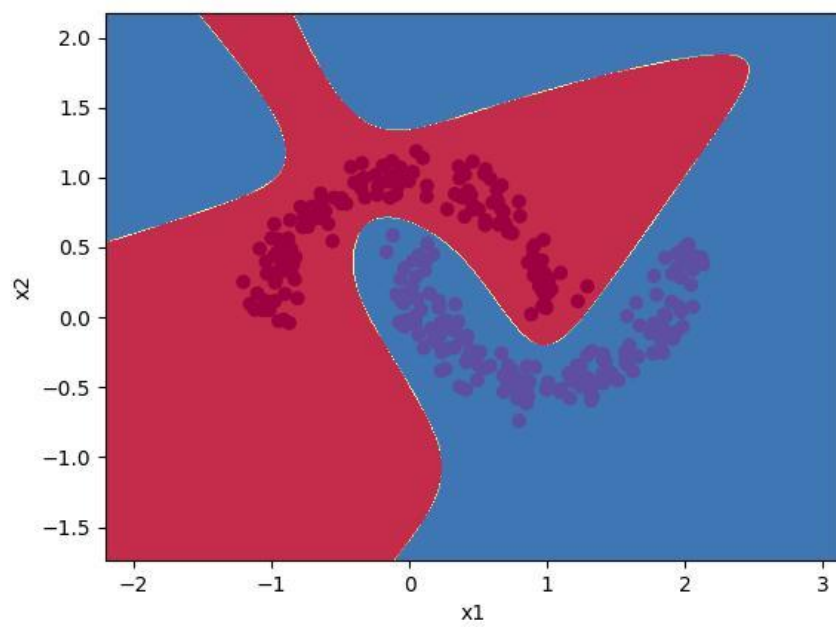
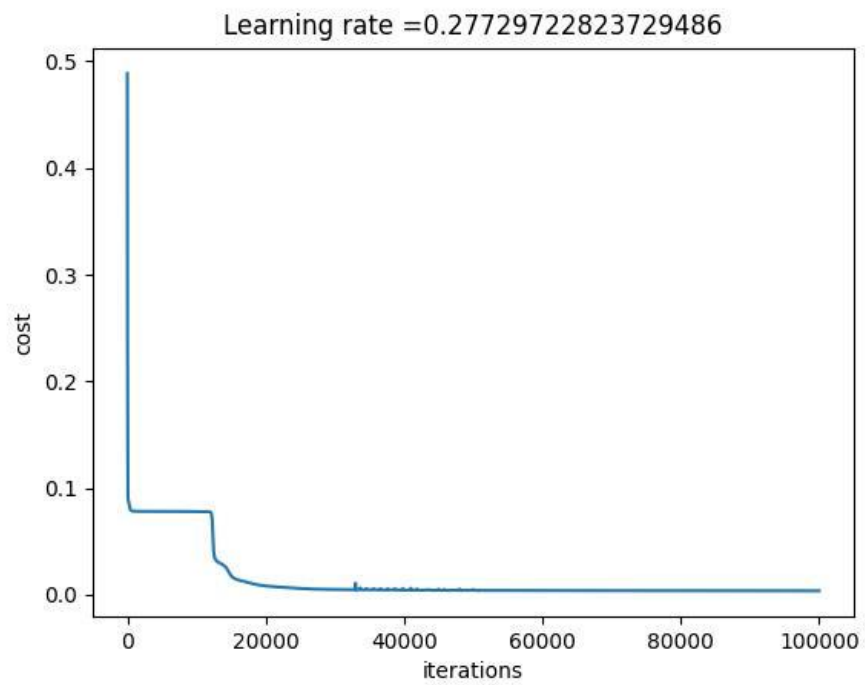
تابع فعالیت ها سیگموئید بوده ، آموزش به صورت batch و نرخ یادگیری ثابت



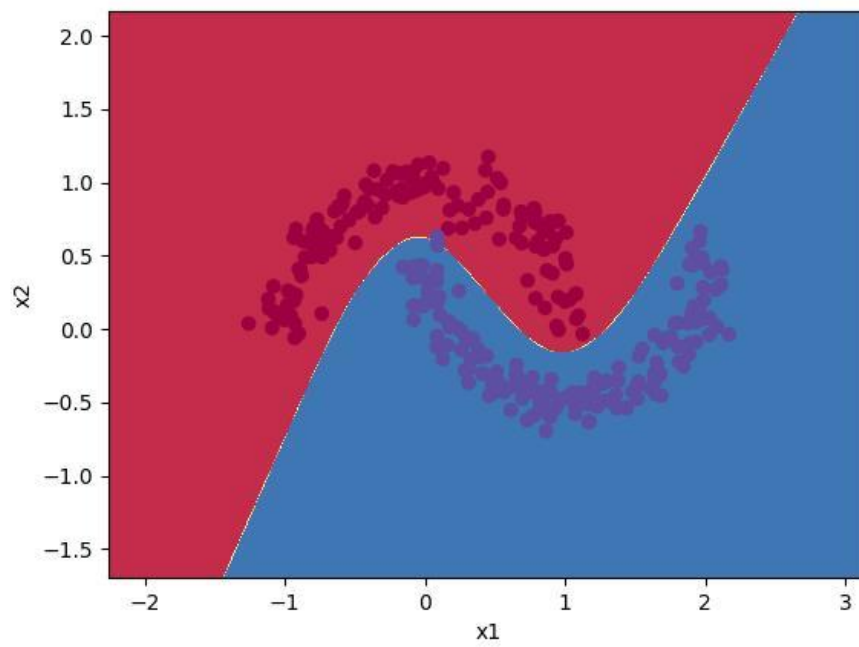
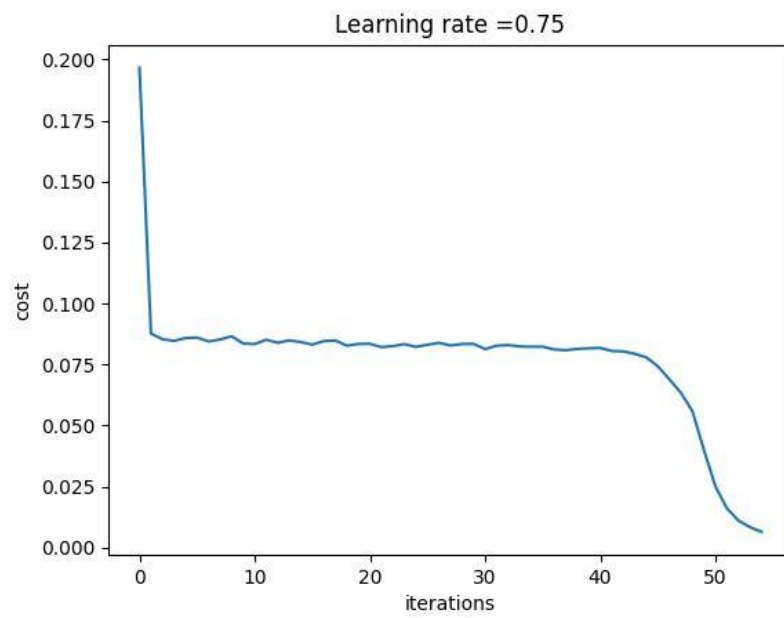
تابع فعالیت tanh ، batch learning و نرخ یادگیری ثابت



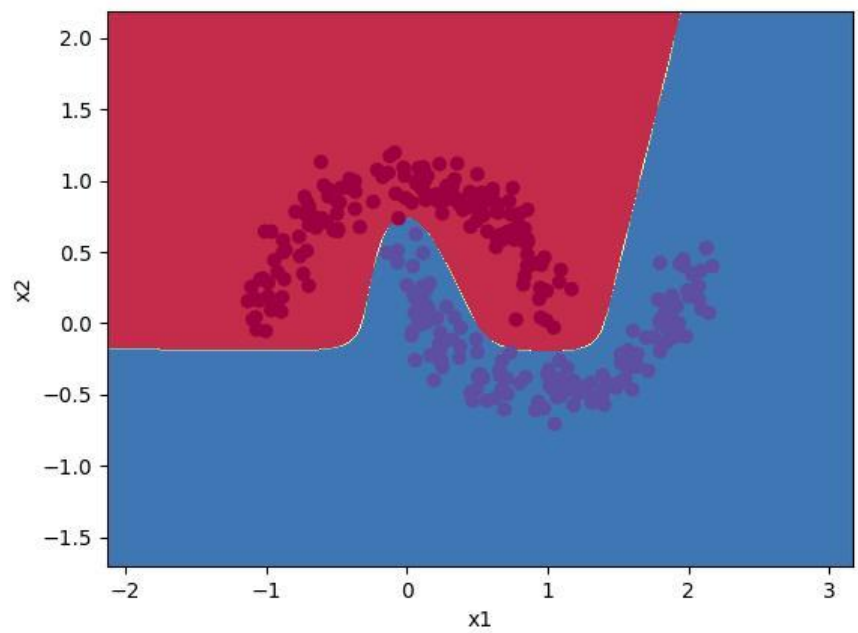
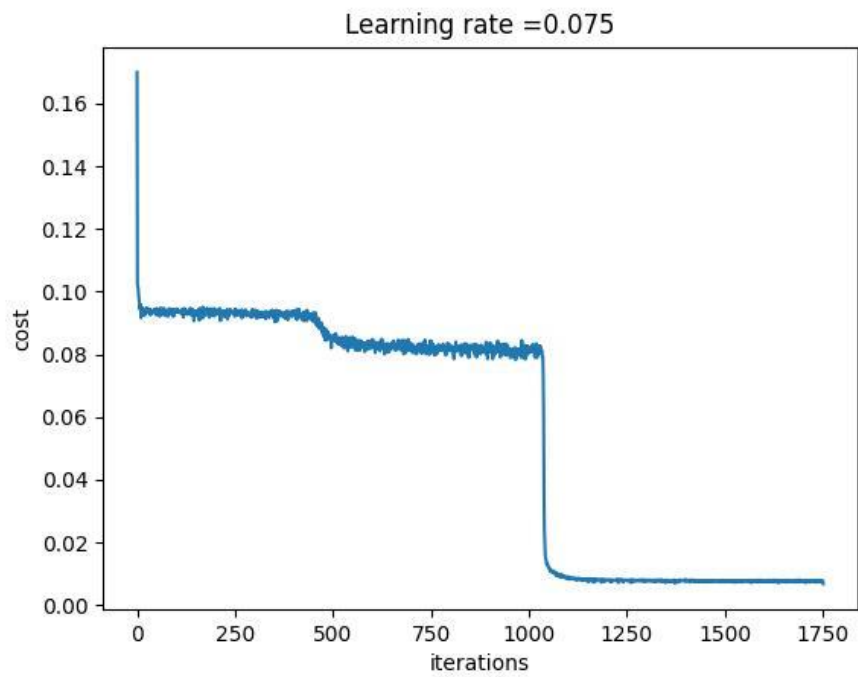
تابع فعاليت sigmoid ، batch learning ، نرخ يادگيري متغير



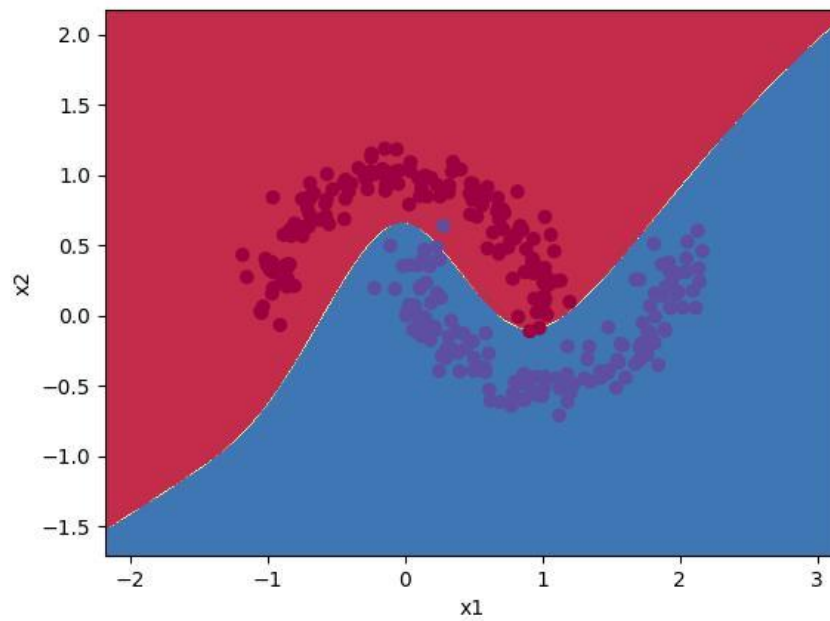
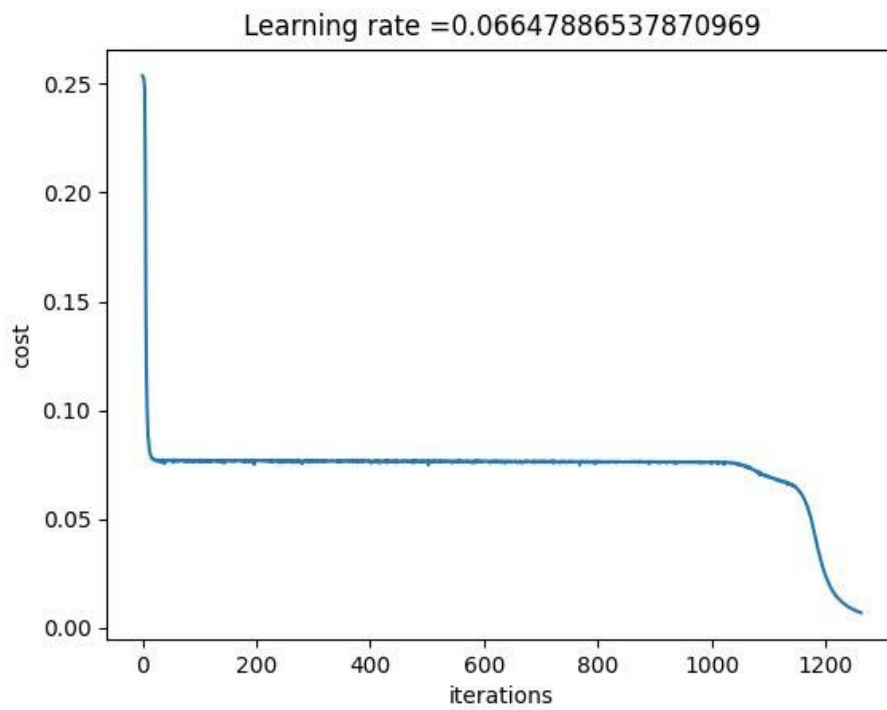
تابع فعاليت \tanh ، batch learning ، نرخ يادگيري متغير



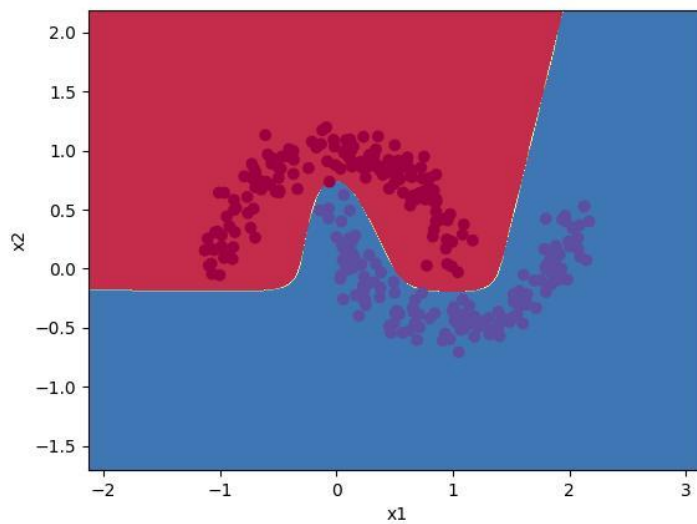
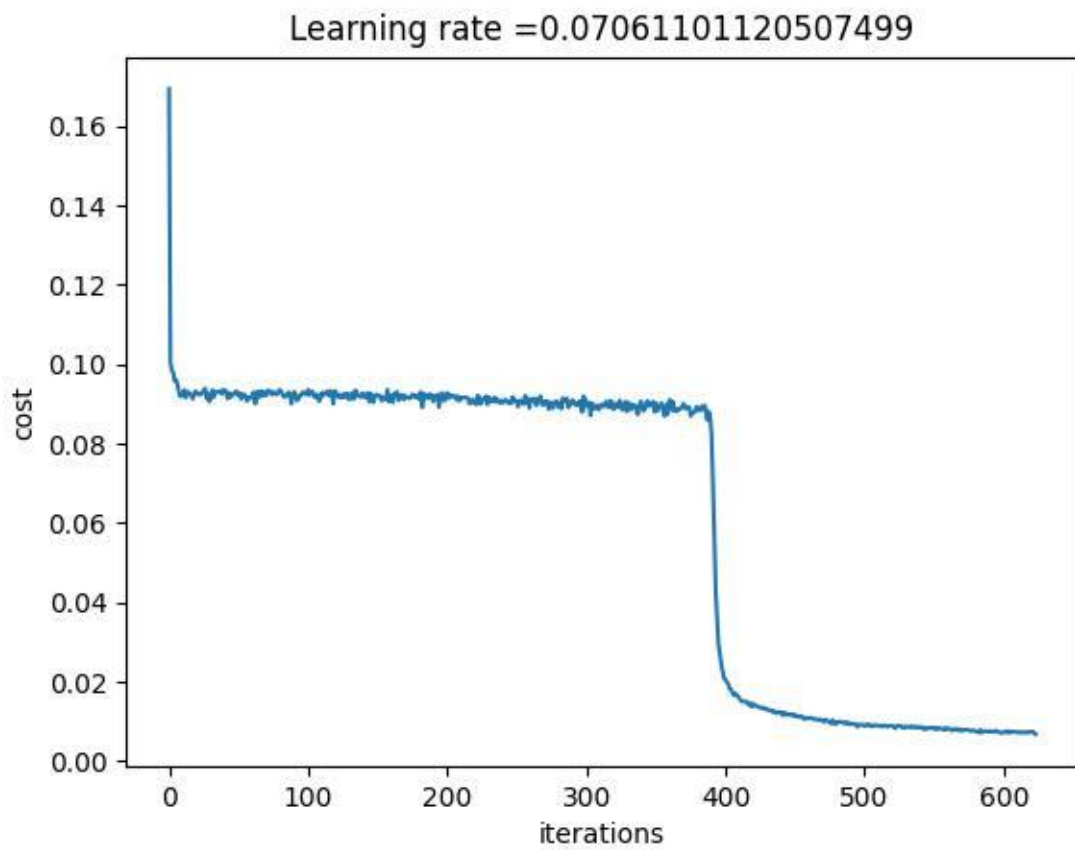
تابع فعاليت sigmoid ، online learning ، نرخ يادگيري ثابت



تابع فعاليت \tanh ، online learning ، نرخ يادگيري ثابت



تابع فعالیت sigmoid ، online learning ، نرخ یادگیری متغیر

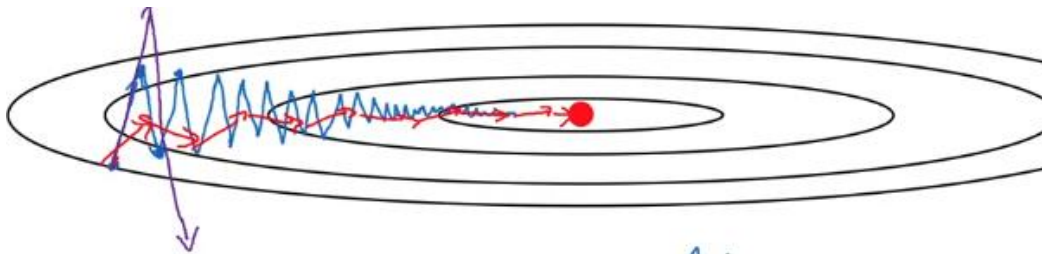


تابع فعاليت \tanh ، online learning ، نرخ يادگيري متغير

به خوبی مشخص است که تقریباً در تمامی روش‌ها توانسته ایم شبکه را به خوبی بر روی داده‌های آموزشی، آموزش دهیم و نمودار نهایی نشان دهنده دقت خوب شبکه بر روی داده‌های تست دارد.

در یادگیری Online مقدار زمان بیشتری نسبت به batch زمان صرف می‌شود چون batch به صورت vectorized شده انجام می‌شود ولی Online هر epoch مقدار زمان بیشتری نیاز دارد.

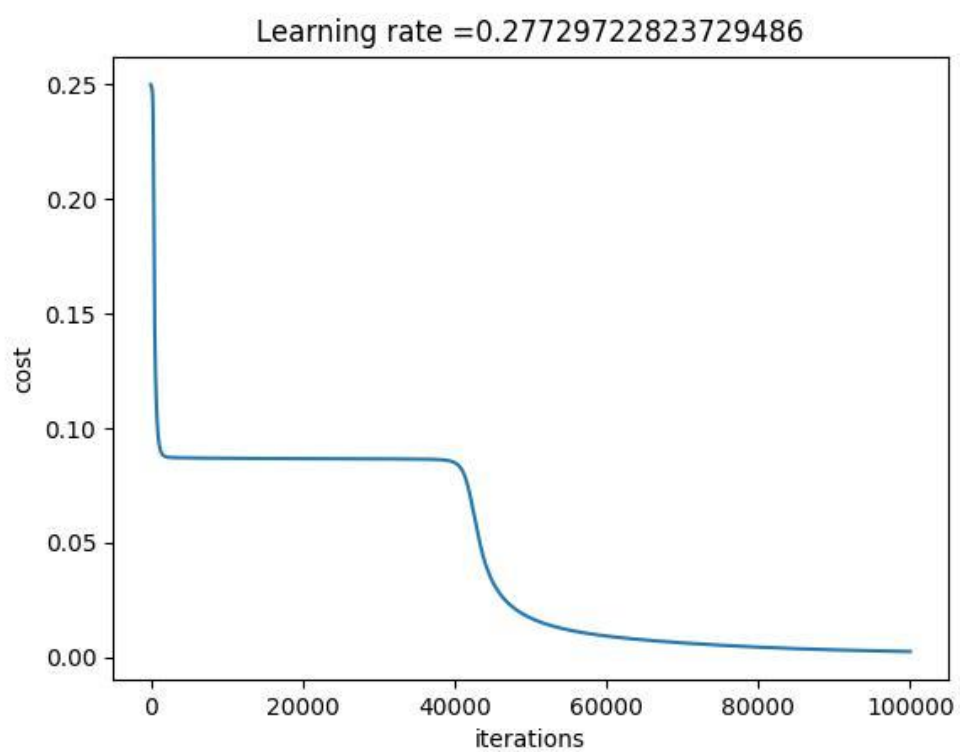
در روش ممنتوم هدف محاسبه یک میانگین وزن دار از گرادیان‌ها و استفاده از آن‌ها برای آپدیت کردن وزن‌های می‌باشد. در ممنتوم هدف آن است که اگر هدف در مرکز کانتور‌ها باشد، آموزش در جهت افقی بیشتر و در جهت عمودی کندتر باشد به نحوی که برخلاف تنزل گرادیان معمول oscillate زیادی نداشته باشد.

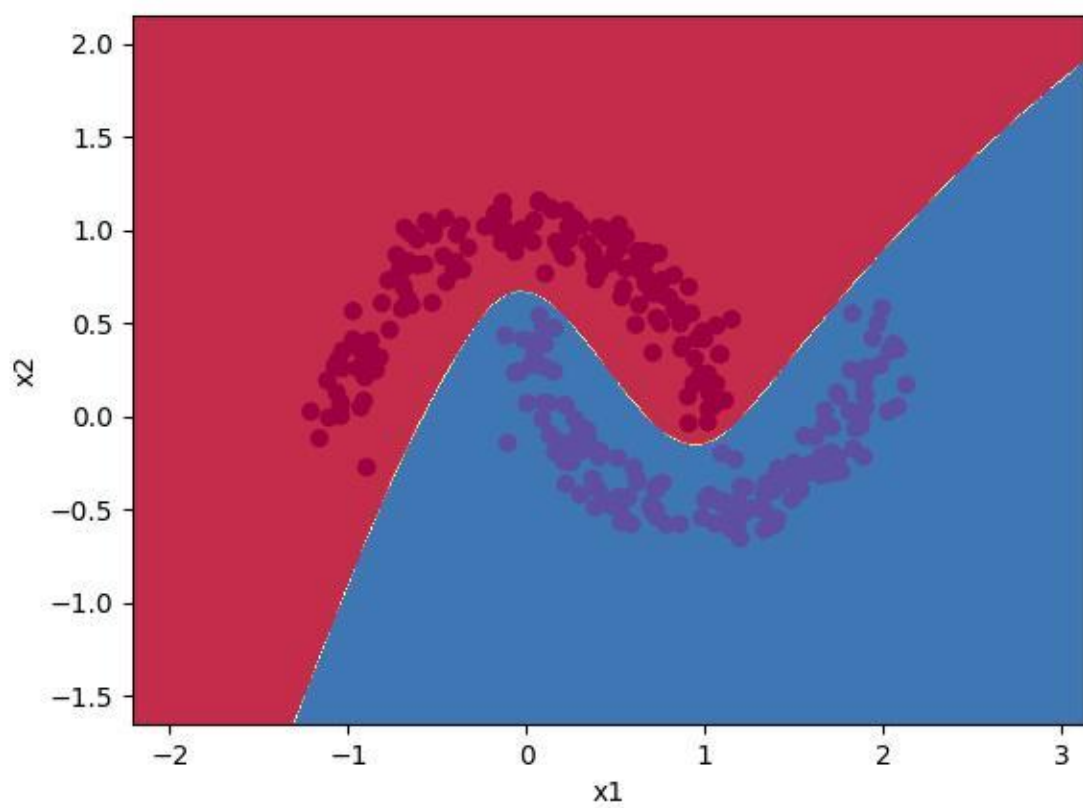


خطوط آبی و بنفش تنزل گرادیان معمول با نرخ یادگیری را نشان می‌دهند و خط قرمز استفاده از ممنتوم را به خوبی نمایش می‌دهد که به صورت یک میانگین وزن دار به سرعت به هدف نزدیک می‌شود.

بهترین شبکه به دست آمده شبکه با تابع فعالیت سیگموئید با نرخ یادگیری متغیر بود.

پس از استفاده از روش ممنتوم، همانطور که در شکل زیر مشخص است سریع‌تر و در epoch‌های کمتری به local optima همگرا شد.





(۲) Adam یا adaptive moment estimation یکی از optimization algorithm هایی است که به خوبی در بازه زیادی از برنامه ها قابل استفاده است . این الگوریتم به زبان ساده به صورت همزمان RMSprop و momentum را در کنار هم قرار داده و از آن ها استفاده میکند. Adam توانسته سرعت روش ممنتوم را با قدرت RMSprop که انطباق گرادیان ها در جهت های مختلف است ترکیب کند.

این الگوریتم در هنگامی که نرخ های یادگیری پارامتر ها منطبق میکند علاوه بر استفاده از میانگین گشتاور اول ما از میانگین گشتاور دوم گرادیان ها هم استفاده میکند. به این معنی که میانگین حرکت نمایی گرادیان و مربع گرادیان را محاسبه میکند. این کار سبب میشود که نرخ کاهش گرادیان به نحوی کنترل شود که کمترین حرکت زیگزاگی oscillation در هنگام رسیدن به گلوبال مینیمم داشته و در همین حال بتوان گام های به اندازه کافی بزرگ برای عبور از مینیمم های محلی داشت .

$$V_{dw} = \text{beta1} * V_{dw} + (1-\text{beta1})dw$$

$$S_{dw} = \text{beta2} * S_{dw} + (1-\text{beta2}) dw^{**2}$$

$$V_{db} = \text{beta1} * V_{db} + (1-\text{beta1})db$$

$$S_{db} = \text{beta2} * S_{db} + (1-\text{beta2})db^{**2}$$

$$V_{dw_corrected} = V_{dw} / (1-\text{beta1}^{**t})$$

$$S_{dw_corrected} = S_{dw} / (1-\text{beta2}^{**t})$$

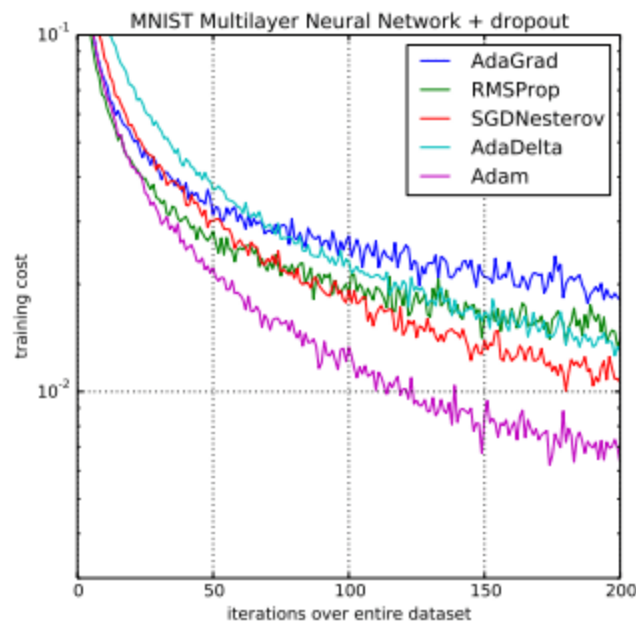
$$V_{db_corrected} = V_{db} / (1-\text{beta1}^{**t})$$

$$S_{db_corrected} = S_{db} / (1-\text{beta2}^{**t})$$

$$W := W - \text{alpha} (V_{dw_corrected})/\text{sqrt} (S_{dw_corrected})+\text{epsilon}$$

$$b := b - \text{alpha} (V_{db_corrected})/\text{sqrt} (S_{db_corrected})+\text{epsilon}$$

Beta1 is the decay rate for the first moment, sum of gradient (aka momentum), commonly set at 0.9. Beta 2 is the decay rate for the second moment, sum of gradient squared, and it is commonly set at 0.999.



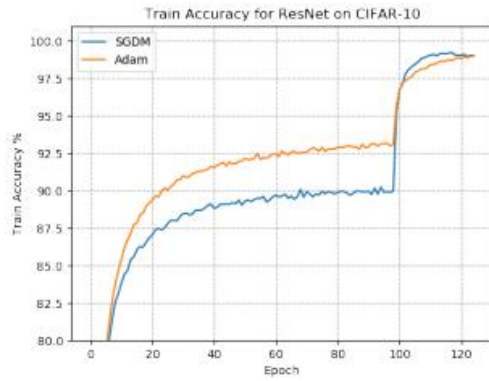
در تصویر بالا عملکرد Adam بر روی دیتاست MNIST مشخص است که بهتر از بقیه optimizer ها مثل RMSprop استفاده میکند.

دلیل برتری Adam نسبت به تنزل گرادیان و ممنتوم از این جهت است که سرعت بالاتری در همگرا شدن دارد و همچنین از نقاط قوت روش ممنتوم نیز استفاده میکند که یکی از علت های سریع تر بودن Adam است .

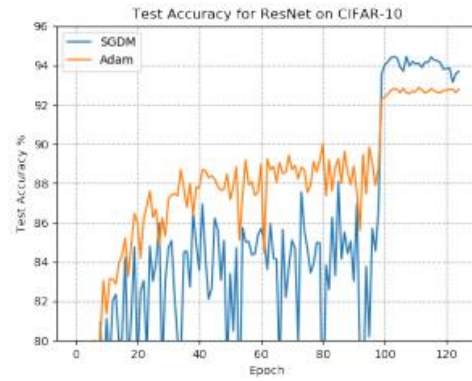
Adam به خوبی مناسب دیتاست هایی با حجم بالا میباشد و همچنین از لحاظ محاسباتی نیز بسیار به صرفه است .

اگرچه Adam به سرعت همگرا میشود ولی میتواند به یک راه حل بهینه همگرا نشود و عملکرد Gradient descent در همگرایی بهتر از Adam باشد زیرا Adam نمیتواند در Generalization مانند SGD عمل کند.

نرخ یادگیری شده توسط Adam در بعضی از شرایط میتواند خیلی کوچک باشد که سبب میشود که نتواند مسیر مناسب برای همگرا شدن پیدا کند. در بعضی از شرایط Adam نرخ یادگیری را بسیار بزرگ میکند که برای یادگیری مناسب نیست.



(a) train



(b) test

<https://machinelearningmastery.com/adam-optimization-algorithm-for-deep-learning>

<https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-cl-۲۳۲۵b۸۹۸momentum-adagrad-rmsprop-adam-f>

<https://www.coursera.org/learn/deep-neural-network/lecture/wVCZ/adam-optimization-۹>

/