



华南理工大学

South China University of Technology

《机器学习》课程实验报告

学 院 软件学院

专 业 软件工程

组 员 吴国斌

学 号 201530613030

邮 箱 854370642@qq.com

指导教师 吴庆耀

提交日期 2017 年 12 月 2 日

1. 实验题目: 逻辑回归、线性分类与随机梯度下降

2. 实验时间: 2017 年 12 月 2 日

3. 报告人: 吴国斌

4. 实验目的:

1. 对比理解梯度下降和随机梯度下降的区别与联系。
2. 对比理解逻辑回归和线性分类的区别与联系。
3. 进一步理解 SVM 的原理并在较大数据上实践。

5. 数据集以及数据分析:

实验使用的是 LIBSVM Data 的中的 a9a 数据, 包含 32561 / 16281(testing)个样本, 每个样本有 123/123 (testing)个属性。请自行下载训练集和验证集。其中 label 表示该记录的类别, +1 为正类, -1 为负类。其中正类有 11687 个样本, 负类有 37155 个样本。

6. 实验步骤:

逻辑回归与随机梯度下降

1. 读取实验训练集和验证集。
2. 逻辑回归模型参数初始化, 可以考虑全零初始化, 随机初始化或者正态分布初始化。
3. 选择Loss函数及对其求导, 过程详见课件ppt。
4. 求得部分样本对Loss函数的梯度 G 。
5. 使用不同的优化方法更新模型参数 (NAG, RMSProp, AdaDelta和Adam)。
6. 选择合适的阈值, 将验证集中计算结果大于阈值的标记为正类, 反之为负类。在验证集上测试并得到不同优化方法的Loss函数值 L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ 和 L_{Adam} 。
7. 重复步骤4-6若干次, 画出 L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ 和 L_{Adam} 随迭代次数的变化图。

线性分类与随机梯度下降

1. 读取实验训练集和验证集。
2. 支持向量机模型参数初始化, 可以考虑全零初始化, 随机初始化或者正态分布初始化。
3. 选择Loss函数及对其求导, 过程详见课件ppt。
4. 求得部分样本对Loss函数的梯度 G 。
5. 使用不同的优化方法更新模型参数 (NAG, RMSProp, AdaDelta和Adam)。
6. 选择合适的阈值, 将验证集中计算结果大于阈值的标记为正类, 反之为负类。在验证集上测试并得到不同优化方法的Loss函数值 L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ 和 L_{Adam} 。
7. 重复步骤4-6若干次, 画出 L_{NAG} , $L_{RMSProp}$, $L_{AdaDelta}$ 和 L_{Adam} 随迭代次数的变化图。

7. 代码内容:

逻辑回归:

```
1 from sklearn import datasets, model_selection, linear_model
2 import numpy as np
3 import jupyter
4 import matplotlib.pyplot as plt
5 import math
6 import random
7
8 X_train, y_train = datasets.load_svmlight_file("a9a_train.txt")
9
10 # turn the csr_matrix into array for further processing
11 x_ = np.array(X_train.toarray(), np.float32).reshape((-1, 123))
12 y_ = np.array(y_train, np.float32).reshape((-1, 1))
13
14 for i in range(y_.shape[0]):
15     if y_[i,0] == -1.0 : y_[i,0] = 0.
16
17 X_ = np.hstack([x_, np.ones((x_.shape[0], 1))])
18
19
20 X_test, y_test = datasets.load_svmlight_file("a9a_test.txt")
21
22 xt_ = np.array(X_test.toarray(), np.float32).reshape((-1, 122))
23 yt_ = np.array(y_test, np.float32).reshape((-1, 1))
24 for i in range(yt_.shape[0]):
25     if yt_[i,0] == -1.0 : yt_[i,0] = 0.
26 # X_ = (x_+1) to fit the constant item
27
28 Xt_ = np.hstack([xt_, np.zeros((xt_.shape[0], 1)), np.ones((xt_.shape[0], 1))])
29
30
31 #  $h_\theta(X) = e^{(\theta * X)} / (1 + e^{(\theta * X)}) = 1 / (1 + e^{(-\theta * X)})$ 
32 def h_theta(Xi, Theta):
33     e_t = math.exp(Xi.dot(Theta.T))
34     return e_t / (1 + e_t)
35
36 #  $L(\theta) = -(1/m) \sum (y_i * \log(h_\theta(X_i)) + (1 - y_i) * \log(1 - h_\theta(X_i)))$ 
37 def compute_loss(X, y, Theta):
38     m = y.shape[0]
39     loss = 0.
40     for i in range(m):
41         loss += (y[i] * math.log(h_theta(X[i, :], Theta))) + ((1 - y[i]) * math.log(1 - h_theta(X[i, :], Theta)))
42     loss /= -m
43     return loss
44
45 #  $\partial L / \partial \theta = (1/m) \sum (h_\theta(X_i) - y_i) * X_i$ 
46 def compute_gradient(X, y, Theta):
47     m = y.shape[0]
48     gradient = np.zeros(Theta.shape)
49     for i in range(m):
50         gradient += (h_theta(X[i, :], Theta) - y[i]) * (X[i, :])
51     gradient /= m
52     return gradient
53
54
55 def train_model_nag(X, y, Theta, learning_rate, gamma, iteration = 10000):
56     test_loss_history = np.zeros((iteration, 1))
57     v = np.zeros(Theta.shape)
58     Theta_gradient = np.zeros(Theta.shape)
59     for iter in range(iteration):
60         index = random.randint(0, y.shape[0]-10)
61         Theta = Theta - gamma * v
62         v = gamma * v - learning_rate * compute_gradient(X[index:index+10, :], y[index:index+10], Theta)
63         Theta = Theta + v
64         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
65     return test_loss_history, Theta
66
67
68 def train_model_rmsprop(X, y, Theta, learning_rate, gamma, epsilon, iteration = 10000):
69     test_loss_history = np.zeros((iteration, 1))
70     G_t = 0.
71     Theta_gradient = np.zeros(Theta.shape)
72     for iter in range(iteration):
73         index = random.randint(0, y.shape[0]-10)
74         Theta_gradient = compute_gradient(X[index:index+10, :], y[index:index+10], Theta)
75         G_t = gamma * G_t + (1 - gamma) * Theta_gradient.dot(Theta_gradient.T)
76         Theta = Theta - (learning_rate / np.sqrt(G_t + epsilon)) * Theta_gradient
77         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
78     return test_loss_history, Theta
79
```

```

80
81 def train_model_adadelta(X, y, Theta, gamma, epsilon, iteration):
82     test_loss_history = np.zeros((iteration, 1))
83     Theta_gradient = np.zeros(Theta.shape)
84     G_t = 0.
85     delta_theta = np.zeros(Theta.shape)
86     delta_t = 0.03
87     for iter in range(iteration):
88         index = random.randint(0, y.shape[0]-10)
89         Theta_gradient = compute_gradient(X[index:index+10,:], y[index:index+10], Theta)
90         G_t = gamma * G_t + (1 - gamma) * Theta_gradient.dot(Theta_gradient.T)
91         delta_theta = - (np.sqrt(delta_t + epsilon) / np.sqrt(G_t + epsilon)) * Theta_gradient
92         Theta = Theta + delta_theta
93         delta_t = gamma * delta_t + (1 - gamma) * (delta_theta.dot(delta_theta.T))
94         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
95     return test_loss_history, Theta
96
97 def train_model_adam(X, y, Theta, learning_rate, beta1, beta2, epsilon, iteration):
98     test_loss_history = np.zeros((iteration, 1))
99     Theta_gradient = np.zeros(Theta.shape)
100     v_t = 0.
101     m_t = np.zeros(Theta.shape)
102     for iter in range(iteration):
103         index = random.randint(0, y.shape[0]-10)
104         Theta_gradient = compute_gradient(X[index:index+10,:], y[index:index+10], Theta)
105         m_t = beta1 * m_t + (1 - beta1) * Theta_gradient
106         v_t = beta2 * v_t + (1 - beta2) * Theta_gradient.dot(Theta_gradient.T)
107         mt_estimate = m_t / (1 - pow(beta1, iter + 1))
108         vt_estimate = v_t / (1 - pow(beta2, iter + 1))
109         Theta = Theta - learning_rate * mt_estimate / (np.sqrt(vt_estimate) + epsilon)
110         test_loss_history[iter] = compute_loss(Xt_, yt_, Theta)[-1:]
111     return test_loss_history, Theta
112
113 iteration = 1000
114
115 be1 = 0.9
116 be2 = 0.999
117 ep = 1e-8
118
119
120
121 t_nag = np.zeros((1, 124))
122 t_rmsprop = np.zeros((1, 124))
123 t_adadelta = np.zeros((1, 124))
124 t_adam = np.zeros((1, 124))
125 # for i in range(t.shape[0]):
126 #     t[i] = [-0.03]
127
128 nag_loss_history, t_nag = train_model_nag(X_, y_, t_nag, 0.005, be1, iteration)
129 rmsprop_loss_history, t_rmsprop = train_model_rmsprop(X_, y_, t_rmsprop, 0.005, be1, ep, iteration)
130 adadelta_loss_history, t_adadelta = train_model_adadelta(X_, y_, t_adadelta, be1, ep, iteration)
131 adam_loss_history, t_adam = train_model_adam(X_, y_, t_adam, 0.005, be1, be2, ep, iteration)
132
133
134 plt.plot(nag_loss_history, 'g', label='NAG')
135 plt.plot(rmsprop_loss_history, 'b', label='RMSProp')
136 plt.plot(adadelta_loss_history, 'r', label='AdaDelta')
137 plt.plot(adam_loss_history, 'y', label='Adam')
138
139
140 plt.legend(loc='upper right')
141
142 plt.ylabel('lost');
143
144 plt.xlabel('iteration count')
145
146 plt.title('loss graph')
147
148 plt.show()

```

线性分类:

```

1 import numpy
2 import random
3 import jupyter
4 import math
5 from sklearn.datasets import load_svmlight_file
6 from sklearn.model_selection import train_test_split
7 from matplotlib import pyplot
8
9 x, y_train = load_svmlight_file("a9a_train.txt")
10 x_train = x.toarray()
11 x, y_test = load_svmlight_file("a9a_test.txt")
12 x_test = x.toarray()
13
14 X_train = numpy.hstack([x_train, numpy.ones((x_train.shape[0], 1))])
15 X_test = numpy.hstack([x_test, numpy.zeros((x_test.shape[0], 1))])
16 X_test = numpy.hstack([X_test, numpy.ones((x_test.shape[0], 1))])
17
18 def compute_grad(x, y, w):
19     gradient = x * (y - x.dot(w.T))
20     return gradient
21
22 def compute_loss(x, y, w, random_i):
23     loss = 0
24     a = len(random_i)
25     for m in range(a):
26         loss += 0.5 * ((y[random_i[m]] - x[random_i[m],:].dot(w.T)) ** 2)
27     return loss/a
28
29 def NAG_train(x, y, x_test, y_test, w, C, lr, gamma, threshold, iteration):
30     vt = numpy.zeros(w.shape)
31     loss_history = []
32     test_loss_history = []
33     random_index = []
34     random_test_index = []
35     for i in range(iteration):
36         random_num = random.randint(0, x.shape[0]-1)
37         random_test_num = random.randint(0, x_test.shape[0]-1)
38         random_index.append(random_num)
39         random_test_index.append(random_test_num)
40     for i in range(iteration):
41         gradient = compute_grad(x[random_index[i],:], y[random_index[i]], w-gamma*vt)
42         vt = gamma*vt - lr*gradient
43         w -= vt
44         loss = compute_loss(x, y, w, random_index)
45         loss_history.append(loss)
46         test_loss_history.append(compute_loss(x_test, y_test, w, random_test_index))
47         if loss < threshold:
48             break
49     return w, loss_history, test_loss_history
50

```

```

51 def RMSProp_train(x, y, x_test, y_test, w, C, lr, gamma, threshold, iteration):
52     Gt = 0
53     loss_history = []
54     test_loss_history = []
55     random_index = []
56     random_test_index = []
57     for i in range(iteration):
58         random_num = random.randint(0, x.shape[0]-1)
59         random_test_num = random.randint(0, x_test.shape[0]-1)
60         random_index.append(random_num)
61         random_test_index.append(random_test_num)
62     for i in range(iteration):
63         gradient = compute_grad(x[random_index[i],:], y[random_index[i]], w)
64         Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient.T)
65         w += lr * gradient / math.sqrt(Gt+1e-8)
66         loss = compute_loss(x, y, w, random_index)
67         loss_history.append(loss)
68         test_loss_history.append(compute_loss(x_test, y_test, w, random_test_index))
69         if loss < threshold :
70             break
71     return w, loss_history, test_loss_history

```

```

73 def AdaDelta_train(x, y, x_test, y_test, w, C, lr, gamma, threshold, iteration):
74     Gt = 0
75     variable_t = 0
76     loss_history = []
77     test_loss_history = []
78     random_index = []
79     random_test_index = []
80     for i in range(iteration):
81         random_num = random.randint(0, x.shape[0]-1)
82         random_test_num = random.randint(0, x_test.shape[0]-1)
83         random_index.append(random_num)
84         random_test_index.append(random_test_num)
85     for i in range(iteration):
86         gradient = compute_grad(x[random_index[i],:], y[random_index[i]], w)
87         Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient.T)
88         variable_w = - math.sqrt(variable_t + 1e-8) * gradient / math.sqrt(Gt + 1e-8)
89         w -= variable_w
90         variable_t = gamma*variable_t + (1-gamma)*variable_w.dot(variable_w.T)
91         loss = compute_loss(x, y, w, random_index)
92         loss_history.append(loss)
93         test_loss_history.append(compute_loss(x_test, y_test, w, random_test_index))
94         if loss < threshold :
95             break
96     return w, loss_history, test_loss_history

```

```

98 def Adam_train(x, y, x_test, y_test, w, C, lr, gamma, threshold, iteration):
99     Gt = 0
100     moment = numpy.zeros((1, x.shape[1]))
101     B = 0.9
102     loss_history = []
103     test_loss_history = []
104     random_index = []
105     random_test_index = []
106     for i in range(iteration):
107         random_num = random.randint(0, x.shape[0]-1)
108         random_test_num = random.randint(0, x_test.shape[0]-1)
109         random_index.append(random_num)
110         random_test_index.append(random_test_num)
111     for i in range(iteration):
112         gradient = compute_grad(x[random_index[i],:], y[random_index[i]], w)
113         moment = B*moment + (1-B)*gradient
114         Gt = gamma*Gt + (1-gamma)*gradient.dot(gradient.T)
115         a = lr * math.sqrt(1 - pow(gamma, iteration)) / (1-pow(B, iteration))
116         w += a * moment / math.sqrt(Gt + 1e-8)
117         loss = compute_loss(x, y, w, random_index)
118         loss_history.append(loss)
119         test_loss_history.append(compute_loss(x_test, y_test, w, random_test_index))
120         if loss < threshold :
121             break
122     return w, loss_history, test_loss_history

```

```

124 iteration = 3000
125 # NAG
126 NAG_w = numpy.zeros((1, X_train.shape[1]))
127 NAG_w, NAG_loss_history, NAG_test_loss_history = NAG_train(X_train, y_train, X_test, y_t
128
129 # RMSProp
130 RMS_w = numpy.zeros((1, X_train.shape[1]))
131 RMS_w, RMS_loss_history, RMS_test_loss_history = RMSProp_train(X_train, y_train, X_test,
132
133 # AdaDelta
134 AdaDelta_w = numpy.zeros((1, X_train.shape[1]))
135 AdaDelta_w, AdaDelta_loss_history, AdaDelta_test_loss_history = AdaDelta_train(X_train,
136
137 #Adam
138 Adam_w = numpy.zeros((1, X_train.shape[1]))
139 Adam_w, Adam_loss_history, Adam_test_loss_history = Adam_train(X_train, y_train, X_test,
140
141
142 pyplot.plot(NAG_test_loss_history, label = 'NAG_validation_loss')
143 pyplot.plot(RMS_test_loss_history, label = 'RMSProp_validation_loss')
144 pyplot.plot(AdaDelta_test_loss_history, label = 'AdaDelta_validation_loss')
145 pyplot.plot(Adam_test_loss_history, label = 'Adam_validation_loss')
146 pyplot.legend(loc='upper right')
147 pyplot.ylabel('loss')
148 pyplot.xlabel('iteration')
149 pyplot.title('graph')
150 pyplot.show()

```

逻辑回归:

8. 模型参数的初始化方法:

模型参数的初始化方法采用的是全零初始化。

9.选择的 loss 函数及其导数:

loss function:

$$L(\theta) = -\frac{1}{m} \sum_{i=0}^n (y_i \log(h_{\theta}(X_i)) + (1 - y_i) \log(1 - \log(h_{\theta}(X_i))))$$

$$\text{其中, } h_{\theta}(X) = \frac{1}{1 + e^{-\theta^T \cdot X}}$$

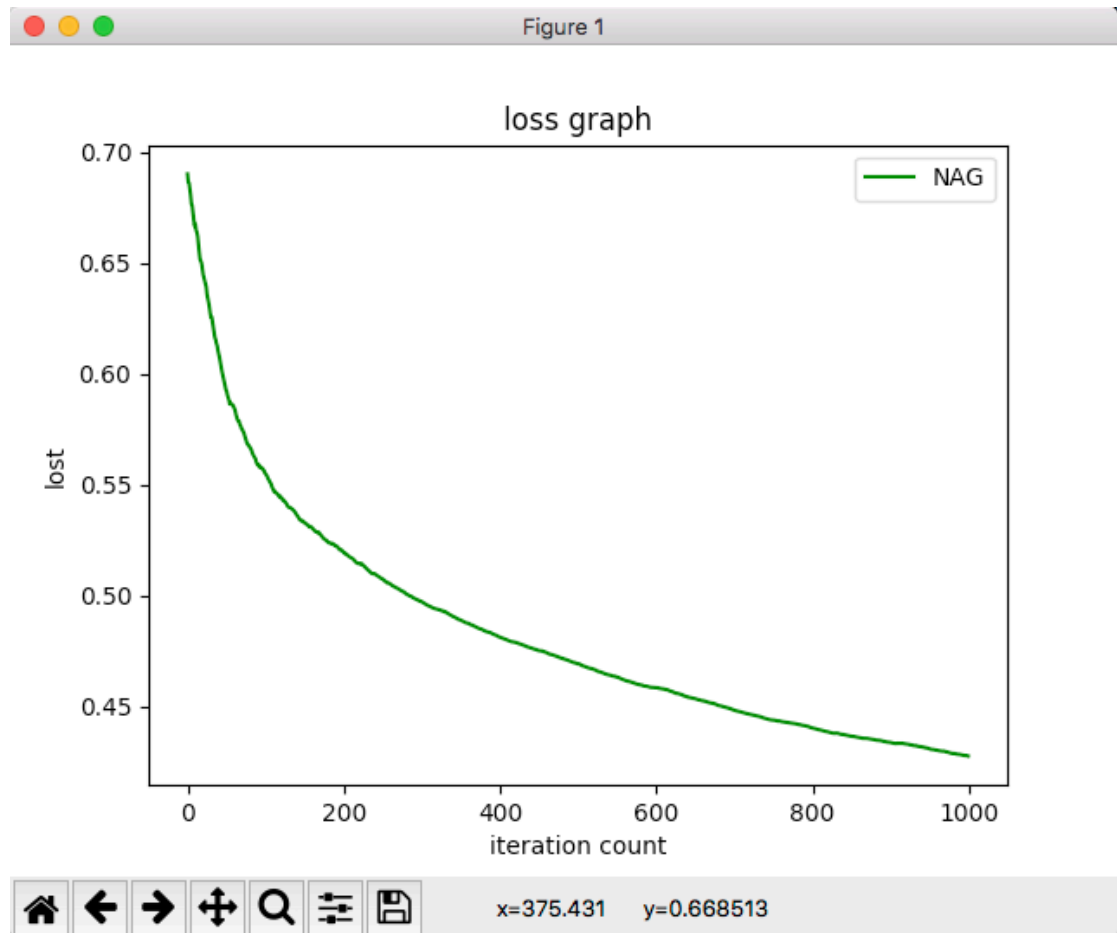
$$\text{gradient: } \frac{\partial L}{\partial \theta} = \frac{1}{n} \sum_{i=0}^n X_i (h_{\theta}(X) - y_i)$$

10.实验结果和曲线图:

NAG:

超参数选择: $\eta=0.005$ $\gamma=0.9$ epoch = 1000

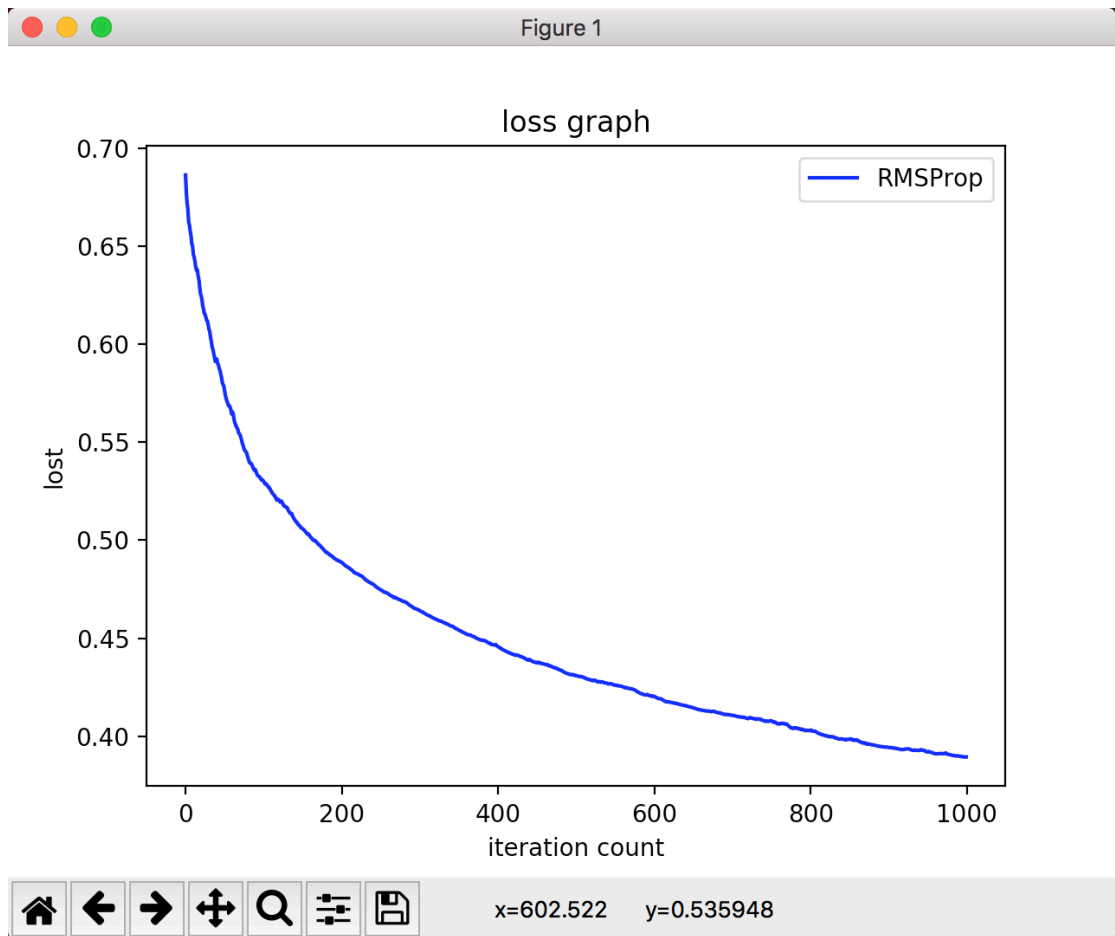
预测结果 (最佳结果):



RMSProp:

超参数选择: $\eta=0.005$ $\gamma = 0.9$ $\varepsilon = 1e-8$ epoch = 1000

预测结果 (最佳结果):



AdaDelta:

超参数选择: $\gamma = 0.9$ $\varepsilon = 1e-8$ epoch=1000

预测结果 (最佳结果):

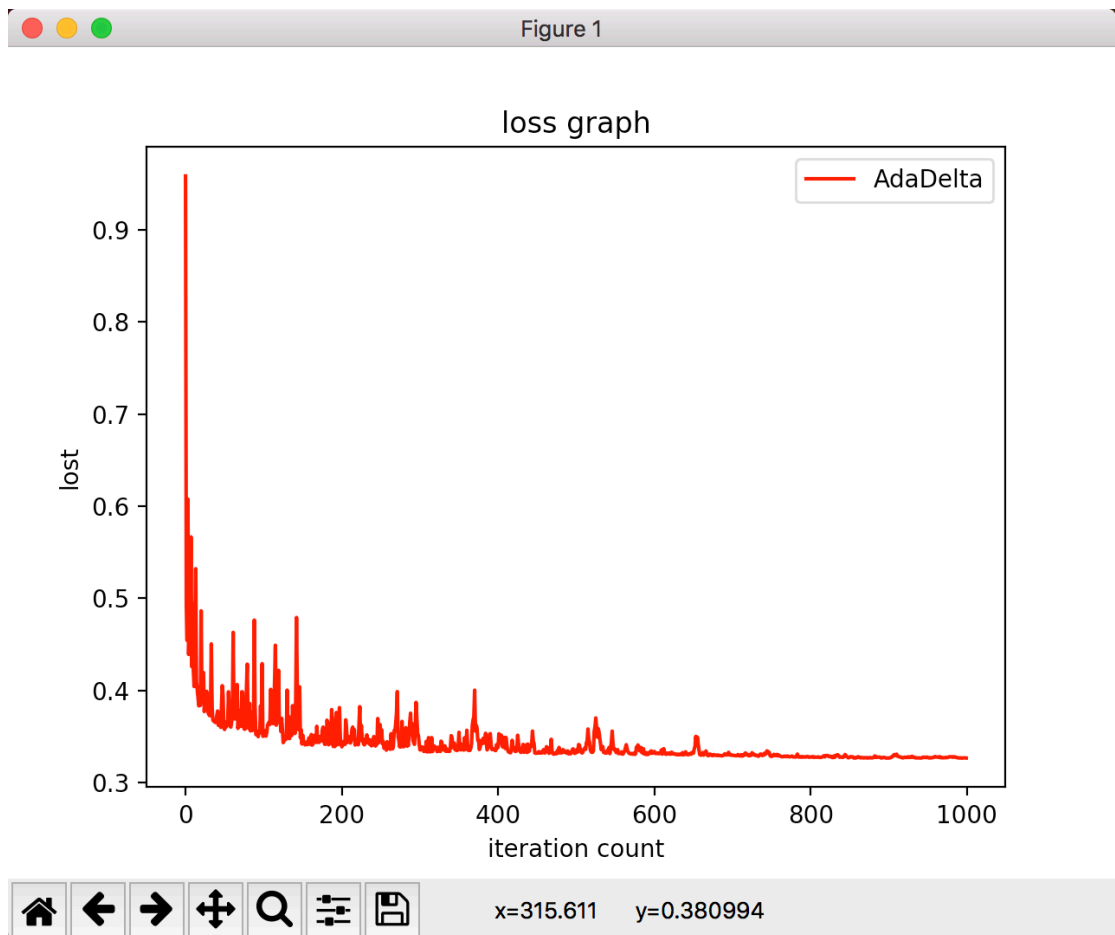
loss 曲线图:

Adam:

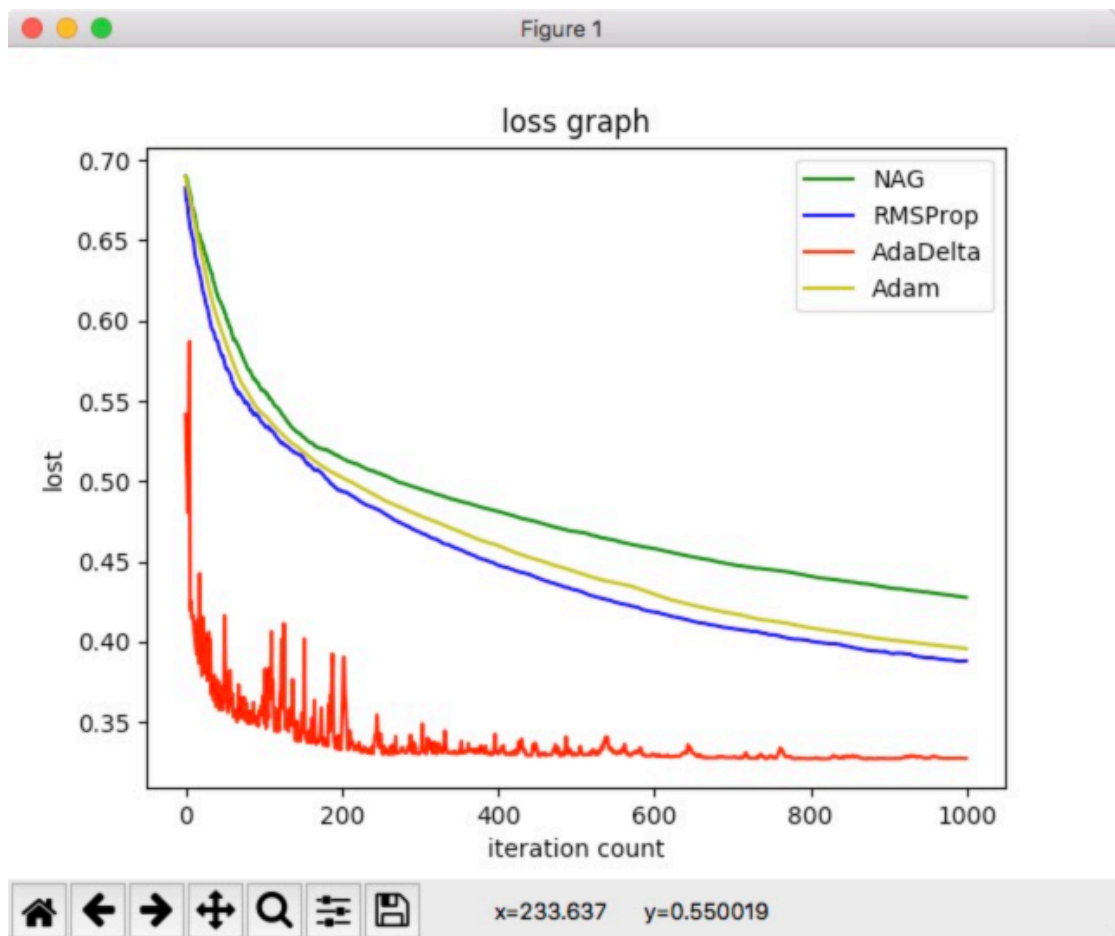
超参数选择: $\eta=0.005$ $\beta_1 = 0.9$ $\beta_2 = 0.999$

$\epsilon = 1e-8$ epoch=1000

预测结果 (最佳结果):



loss 曲线图：



11.实验结果分析:

AdaDelta 较快的达到了局部最优解，但是震荡明显；NAG 在收敛速度和收敛幅度上都较小，但是曲线较为平滑；RMSProp 和 Adam 基本重合，收敛速度比 NAG 快，但是比 AdaDelta 慢。

总体而言，四种方法都可以较为迅速的收敛到局部最优解，比基本的梯度下降法更加优秀。

线性分类:

8. 模型参数的初始化方法:

模型参数的初始化方法采用的是全零初始化。

9.选择的 loss 函数及其导数:

loss function:

$$L(\theta) = \frac{1}{2n} \sum_{i=0}^n (y_i - h_{\theta}(X_i))^2$$

其中, $h_{\theta}(X) = \sum_{i=0}^n \theta_i X_i$

gradient:

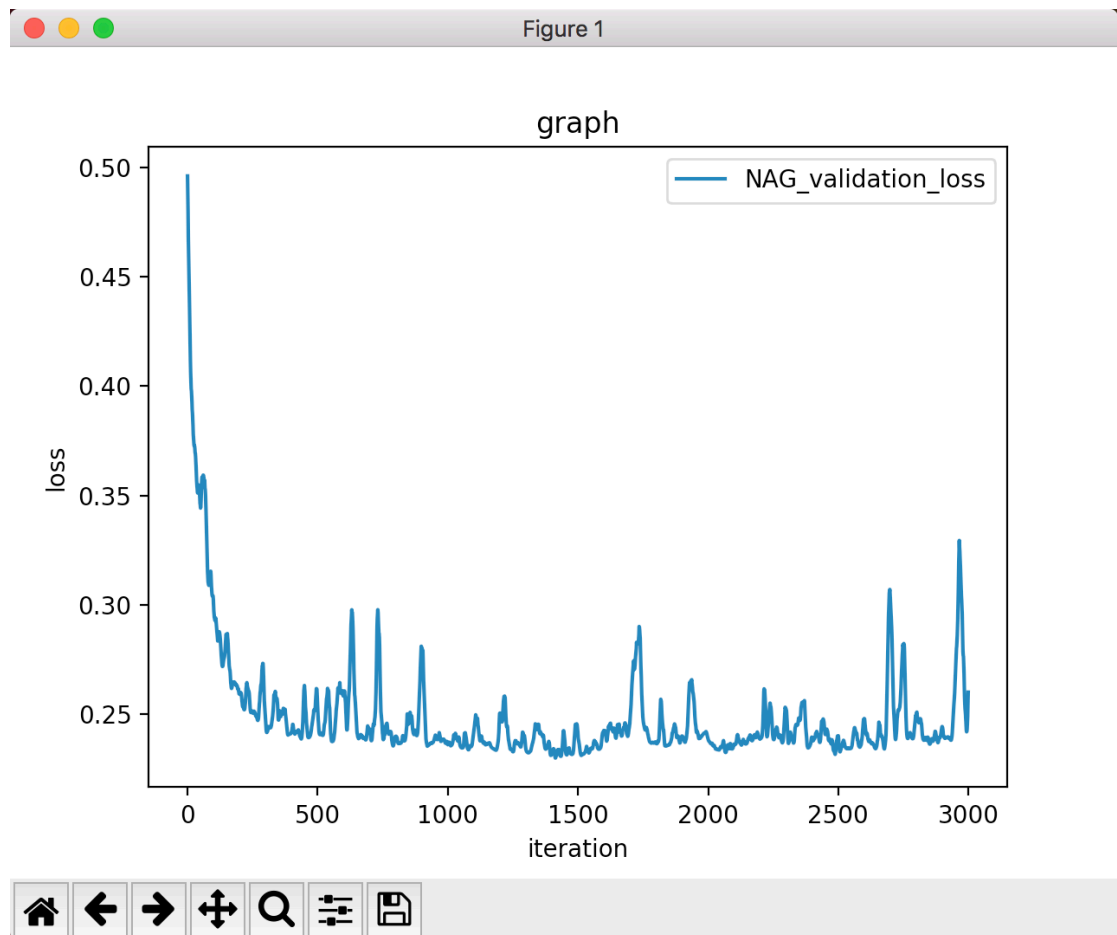
$$\frac{\partial L}{\partial \theta} = \frac{1}{n} \sum_{i=0}^n (y_i - h_{\theta}(X_i)) \cdot X_i$$

10.实验结果和曲线图:

NAG:

超参数选择: $\eta=0.001$ $\gamma=0.9$ epoch = 3000

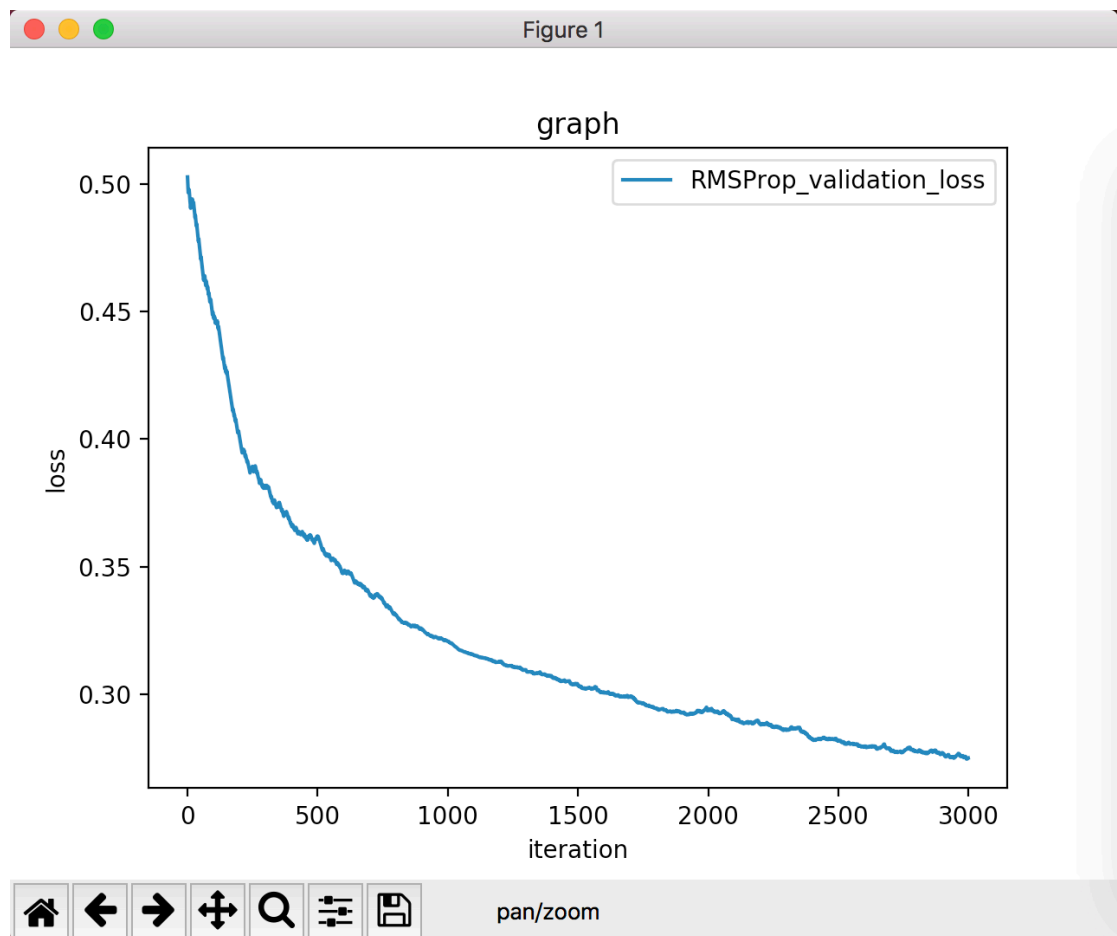
预测结果 (最佳结果):



RMSProp:

超参数选择: $\eta=0.001$ $\gamma=0.9$ $\varepsilon=1e-8$ epoch = 3000

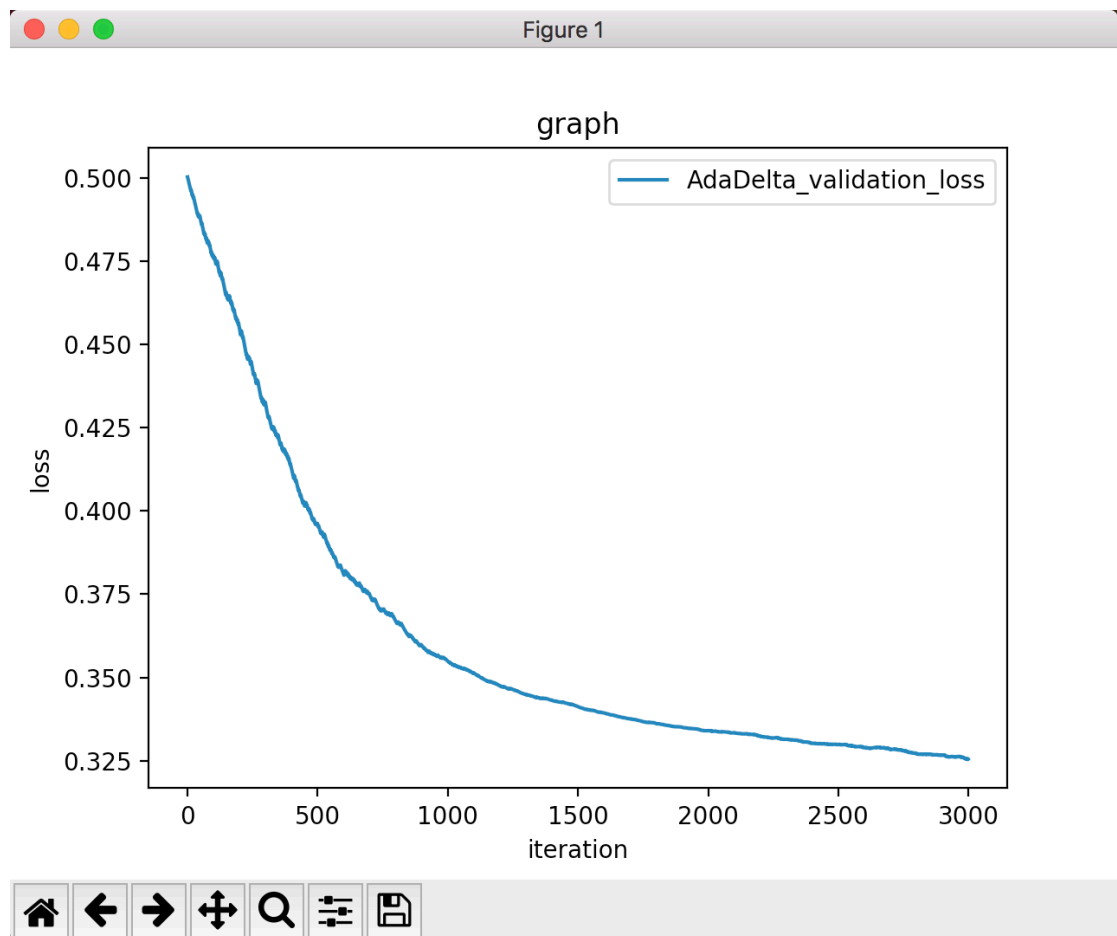
预测结果 (最佳结果):



AdaDelta:

超参数选择: $\gamma = 0.9$ $\varepsilon = 1e-8$ epoch=3000

预测结果 (最佳结果):

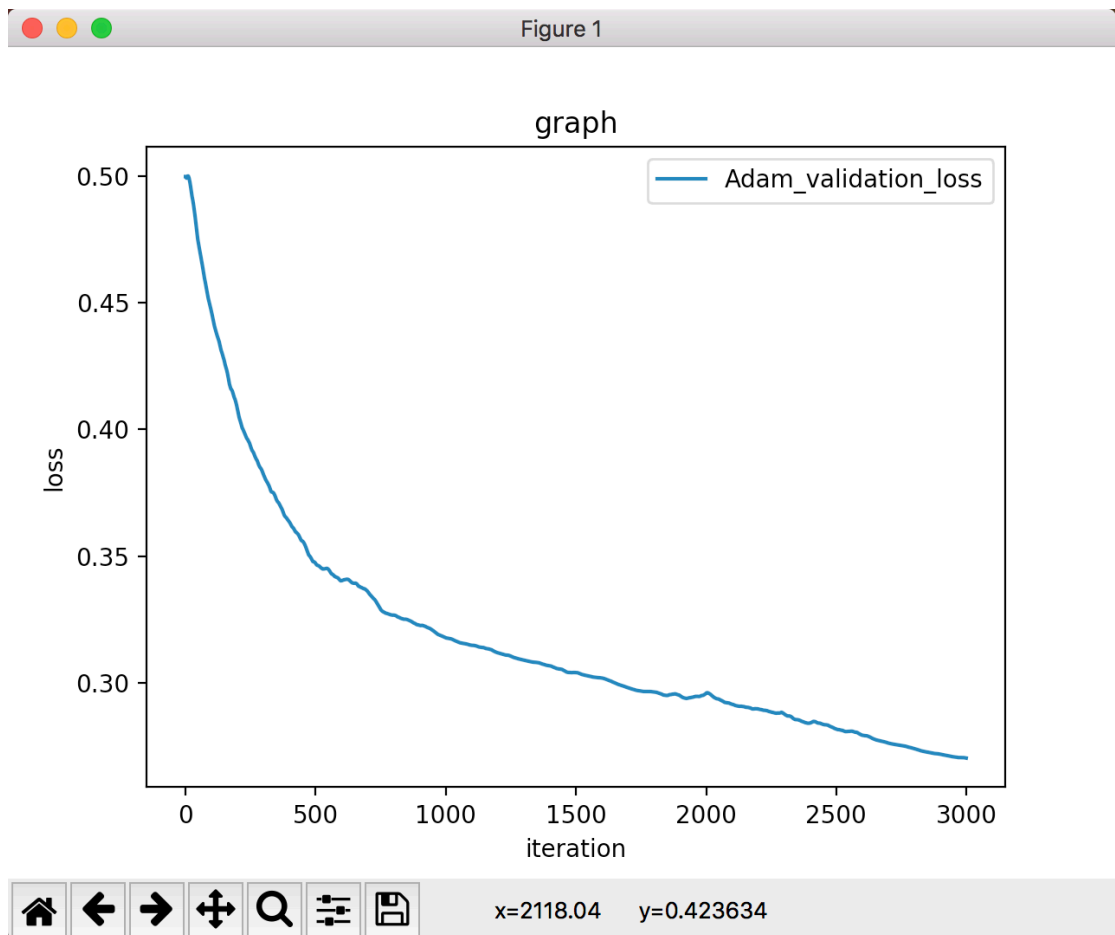


Adam:

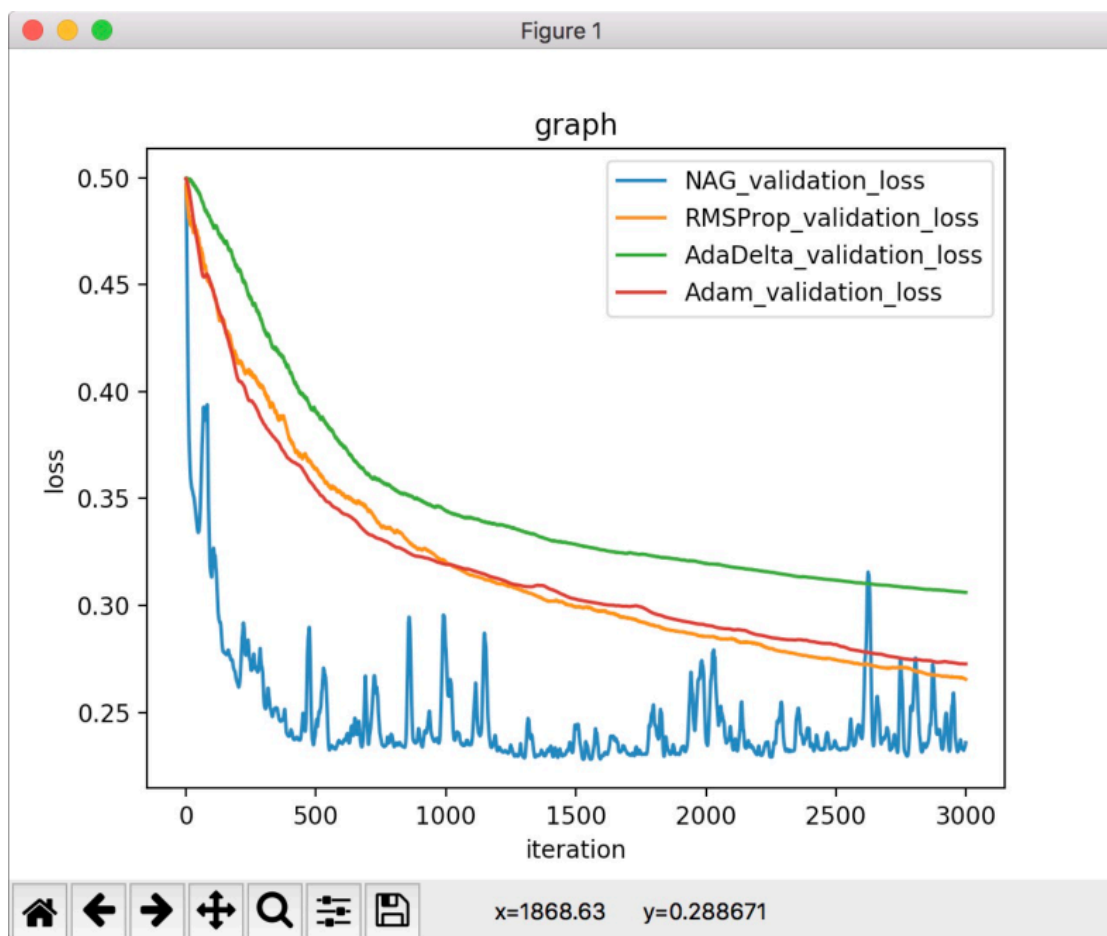
超参数选择: $\eta=0.005$ $\beta_1 = 0.9$ $\beta_2 = 0.999$

$\epsilon = 1e-8$ epoch=3000

预测结果 (最佳结果):



loss 曲线:



11.实验结果分析:

NAG 收敛速度最快，但是震荡明显；； RMSProp 与 Adam 的收敛速度几乎持平，震荡不明显； AdaDelta 收敛较慢但是曲线比较平滑。相比基本的随机梯度下降算法，这四种算法在不同程度上做了改进，更加优秀。

12.对比逻辑回归和线性分类的异同点:

同:

都属于分类问题，都用于预测。

异:

找最优超平面的方法不同，， 形象点说， logistic 模型找的

那个超平面，是尽量让所有点都远离它，而 SVM 线性分类寻找的那个超平面，是只让最靠近中间分割线的那些点尽量远离，即只用到那些“支持向量”的样本。

逻辑回归只可以处理线性可分情况；SVM 则二者皆可。

13.实验总结：

本次实验当中，我学习到了很多 SGD 在实践运用的经验，将课程中学习到的知识运用在实际问题上。但是由于自身对于知识把握懂得程度不高，在实现的过程中遇到了诸如无法正确实现 SGD 优化算法，调参不够灵活的问题，在总结反思之后解决了问题并顺利完成了实验。、

在对模型的训练过程中，我体会到了灵活调参的重要性。在一开始，因为超参数 C, learning rate 等设置得不合理，导致 loss 图像与预期相差甚远，模型参数无法收敛或者收敛过慢，跑数据集时间过长等问题，导致无法拟合数据；或者是因为迭代次数过少，参数还未收敛便停止了训练。在进行数次的不同的调参后，模型往预期方向改变，我也从中学得一些经验。