



本科生实验报告

学生姓名： 丁晓琪

学生学号： 22336057

专业名称： 计科

一：实验任务

二：实验过程

1: 在不同的核函数的SVM二分类器

(1).SVM的基本理论:

(2).SVM的代码实现:

(3).不同核函数的比较分析:

2: 实现采用hinge loss的线性分类模型:

(1). 基于hinge loss的线性分类理论:

(2). 代码实现:

(3). 采用 hinge loss 的线性分类模型和 SVM 模型之间的关系:

3: 采用 hinge loss 线性分类模型和 cross-entropy loss 线性分类模型比较:

(1). cross-entropy loss交叉熵线性模型理论 (逻辑回归) :

(2). 代码实现:

(3). 和 hinge loss 线性分类模型比较:

三：实验结果

1. 不同超参数设置:

2. 所有实验结果:

一：实验任务

1. 考虑两种不同的核函数: i) 线性核函数; ii) 高斯核函数
2. 可以直接调用现成 SVM 软件包来实现
3. 手动实现采用 hinge loss 和 cross-entropy loss 的线性分类模型, 并比较它们的优劣

二：实验过程

1: 在不同的核函数的SVM二分类器

(1).SVM的基本理论:

- 线性最大余量分类器：

- 分界线： $w^T x + b = 0$

- 余量定义：所有样本点距离分类线的最小距离 $\min_l \frac{y^{(l)} \cdot (w^T x^{(l)} + b)}{\|w\|}$

- 训练目标：找到 w^* 和 b^* 使得余量最大，这样基于训练数据训练出的分界线能够更好对未知数据进行判断

- 问题转化：由于求解 $w^*, b^* = \operatorname{argmax}_{w,b} \{ \min_l \frac{y^{(l)} \cdot (w^T x^{(l)} + b)}{\|w\|} \}$ ，既有求max也要求min，可以将问题转化为

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{约束: } y^{(l)} \cdot (w^T x^{(l)} + b) \geq 1, \text{ 对所有训练样本点 } l$$

- 将带约束的优化问题转化为等效对偶问题：注意a为一个训练样本总量维度的向量

$$\max_a g(a)$$

$$g(a) = \sum_{l=1}^N a_l - \frac{1}{2} \sum_{l=1}^N \sum_{j=1}^N a_l a_j y^{(l)} y^{(j)} x^{(l)T} x^{(j)}$$

$$\text{约束: } a \geq 0, \sum_{l=1}^N a_l y^{(l)} = 0$$

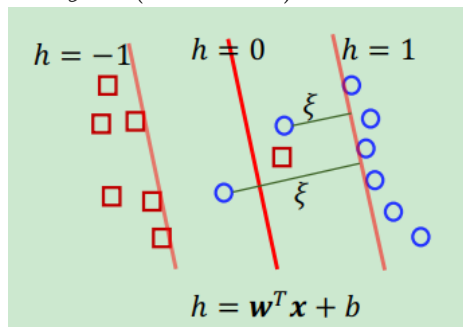
- 当求出上述问题的最优a,可求得最优 w^*

$$w^* = \sum_{l=1}^N a_l y^{(l)} x^{(l)}$$

$$\text{则分类器变成: } \hat{y}(x) = \operatorname{sign}(\sum_{l=1}^N a_l^* y^{(l)} x^{(l)T} x + b^*)$$

- 软间隔最大余量分类器：

- 背景：训练样本可能不能用线性分类器完全区分成两类，可能出现在问题转化步骤出现不满足约束 $y^{(l)} \cdot (w^T x^{(l)} + b) \geq 1$ 的样本点，使得永远无法求得最优解



- 解决：将约束放宽给予松弛量 ξ , ($\xi \geq 0$)

$$\text{约束: } y^{(n)}(w^T x^{(n)} + b) \geq 1 - \xi$$

$$\text{转化问题: } \min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{l=1}^n \xi_n$$

- 等效对偶问题:

Using the same method as before, the dual formulation can be derived as

$$\max_{\mathbf{a}} g(\mathbf{a})$$

$$s. t.: a_n \geq 0, a_n \leq C$$

$$\sum_{n=1}^N a_n y^{(n)} = 0$$

$$\text{where } g(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)}$$

When $a_n > C$, it can be shown that $g(\mathbf{a}) = -\infty$

- SVM:

- 背景: 线性分类器不能对所有样本分类, 引入更高维的非线性分类器
- 关键: 需要 $\phi: x \rightarrow \phi(x)$ 将样本的特征值投影到更高维
- 问题转化: 用非线性的 $\phi(x)$ 代替 x , 且 $\phi(x)$ 越高维度越好, 但是计算上会很复杂

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$s. t.: y^{(n)} \cdot (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) + b) \geq 1 - \xi_n,$$

$$\xi_n \geq 0, \quad \text{for } n = 1, 2, \dots, N$$

- 等效对偶问题:

The problem can be solved via its dual form:

$$\max_{\mathbf{a}} g(\mathbf{a})$$

$$s. t.: a_n \geq 0, a_n \leq C$$

$$\sum_{n=1}^N a_n y^{(n)} = 0$$

$$\text{where } g(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m y^{(n)} y^{(m)} \phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(m)})$$

$$= \mathbf{1}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{M} \mathbf{a}$$

- 分类器:

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* y^{(n)} \phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}) + b^* \right)$$

- 核函数:

- 背景: 即使SVM用等效对偶问题求解时, 需要求解的最优化变量 \mathbf{a}^* 的维度只和样本的数量有关, 但是 $g(\mathbf{a})$ 中仍然包含着需要复杂计算的 $\phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(n)})$
- 解决: 为了优化计算, 可以将 $\phi(\mathbf{x}^{(n)})^T \phi(\mathbf{x}^{(n)})$ 的向量矩阵运算转化为更简单的运算 $k(x, x')$, 该函数只需要输入 x, x' 就可以以更简便的方式计算 $\phi(\mathbf{x})^T \phi(\mathbf{x}')$

(2).SVM的代码实现:

- 初始化:

将训练数据和测试数据从csv文件中提取出来, 用 `StandardScaler()` 对训练数据的特征值和测试数据的特征值标准化

```
1 def __init__(self, filename_train, filename_test):
2     # Parameters:
3     # filename_train (str): 训练数据的csv文件名字。
```

```

4      # filename_test (str): 测试数据的csv文件名字。
5      # 实例变量:
6      # - self.X_train_std: 标准化后的训练数据特征值。
7      # - self.Y_train: 训练数据的标签。
8      # - self.X_test_std: 标准化后的测试数据特征值。
9      # - self.Y_test: 测试数据的标签。
10     X_train,self.Y_train=loadCSVfile1(filename_train)
11     X_test,self.Y_test=loadCSVfile1(filename_test)
12     sc=StandardScaler()
13     sc.fit(X_train)
14     self.X_train_std=sc.transform(X_train)
15     self.X_test_std=sc.transform(X_test)

```

- 用线性核函数在训练数据上训练模型，记录训练时间，且在测试数据上测试模型准确率
直接调用python的SVC类，且指定内核为liner， `svm.fit()` 指定训练数据； `svm.predict()` 用训练好的模型预测标签，返回所有预测的标签值； `svm.score()` 比较模型对测试集的预测结果与实际的标签，并返回一个准确率

```

1      def linear_train(self):
2          #线性核svm分类器
3          #输出: 模型训练时长，在测试样本上的预测正确率和正确预测数
4          start_time = time.time()
5          svm=SVC(kernel='linear',C=1.0,random_state=1)
6          svm.fit(self.X_train_std,self.Y_train)
7          end_time = time.time()
8          training_duration = end_time - start_time
9          print(f"线性核SVM训练时长: {training_duration:.4f} 秒")
10
11         y_pred=svm.predict(self.X_test_std)
12
13         print('linear Misclassified samples: %d' % (self.Y_test !=
14         y_pred).sum())
15         print('linear Accuracy: %.6f' % svm.score(self.X_test_std,
16         self.Y_test))

```

- 用高斯核函数在训练数据上训练模型，记录训练时间，在测试数据上测试模型准确率，指定内核类型为rbf

```

1      def gauss_train(self):
2          #高斯核svm分类器
3          #输出: 模型训练时长，在测试样本上的预测正确率和正确预测数
4          start_time = time.time()
5          svm=SVC(kernel='rbf',C=1.0,random_state=2)
6          svm.fit(self.X_train_std,self.Y_train)
7          end_time = time.time()
8          training_duration = end_time - start_time
9          print(f"高斯核SVM训练时长: {training_duration:.4f} 秒")
10
11         y_pred=svm.predict(self.X_test_std)
12
13         print('gauss Misclassified samples: %d' % (self.Y_test !=
14         y_pred).sum())
15         print('gauss Accuracy: %.6f' % svm.score(self.X_test_std,
16         self.Y_test))

```

(3).不同核函数的比较分析:

- 线性核: $K(x^i x^j) = x^{iT} x^j$, 和软约束最大余量线性分类器完全一样, 没有对数据往更高维投射, 保留了原始特征空间的结构, 计算简单
- 高斯核: $k(x, x') = \exp(-\frac{1}{2\sigma^2} \|x - x'\|^2)$, 同时将特征空间投射为

$$\phi(x) = e^{-x^2/2\sigma^2} \left[1, \sqrt{\frac{1}{1!\sigma^2}} x, \sqrt{\frac{1}{2!\sigma^4}} x^2, \sqrt{\frac{1}{3!\sigma^6}} x^3, \dots \right]^T$$

- 用线性核训练的结果 (改变多次随机种子):

```
线性核SVM训练时长: 0.5570 秒  
linear Misclassified samples: 1  
linear Accuracy: 1.00  
线性核SVM训练时长: 0.5319 秒  
linear Misclassified samples: 1  
linear Accuracy: 0.999527  
线性核SVM训练时长: 0.5947 秒  
linear Misclassified samples: 1  
linear Accuracy: 0.999527
```

- 用高斯核训练的结果:

```
高斯核SVM训练时长: 2.7388 秒  
gauss Misclassified samples: 8  
gauss Accuracy: 1.00  
高斯核SVM训练时长: 2.7062 秒  
gauss Misclassified samples: 8  
gauss Accuracy: 0.996217  
高斯核SVM训练时长: 2.7532 秒  
gauss Misclassified samples: 8  
gauss Accuracy: 0.996217
```

- 比较分析:
 - 高斯核比线性核适合数据非线性可分的情况, 高斯核可以映射到高维特征空间, 找到超平面分类
 - 高斯核的特征空间是无限维的, 非线性的分类能力非常强大
 - 线性核的运算比高斯核的运算更为简单, 适合用在数据线性可分的场景
 - 实验中多次训练结果显示线性核在测试集的正确率比高斯核高一点, 这说明实验的数据是线性可分的。而且线性核的模型训练时间比高斯核时间短, 表现了线性核的计算更为简单

2: 实现采用hinge loss的线性分类模型:

(1). 基于hinge loss的线性分类理论:

- hinge loss定义:

$$L(w, b) = C \sum_{n=1}^N E_{SV}(y^{(n)} h^{(n)}) + \frac{1}{2} \|w\|^2$$

$$E_{SV}(z) = \max(0, 1 - z)$$

$$h(x) = wx + b$$

- 基于hinge loss的线性分类器: $y = wx + b$ $y \in 1, -1$,
目标: 找到 w, b 使得 $L(w, b)$ 最小。当 $y=1$ 时, 想要正确分类, 就需要 $wx + b > 0$, 且要离0远一点, 如果 $wx + b > 1$, 则没有任何惩罚, 否则惩罚为 $1 - wx - b$
- 训练方法: 梯度下降

$$w_i = w_i - \alpha \frac{\sigma L(w, b)}{\sigma w_i}$$

$$E_{SV}(h(x), y) = \max(0, 1 - y h(x))$$

$$L(w, b) = \sum_{i=1}^N E_{SV}(h(x^{(i)}), \hat{y}^{(i)}) + \frac{1}{2} \|w\|^2$$

$$\frac{\partial L(w, b)}{\partial w_i} = \sum_{n=1}^N \frac{\partial E_{SV}(h(x^{(n)}), \hat{y}^{(n)})}{\partial w_i} + \frac{\partial \frac{1}{2} \|w\|^2}{\partial w_i} = \sum_{n=1}^N -\hat{y}^{(n)} x_i^{(n)} + w_i$$

$$\frac{\partial E_{SV}(h(x^{(n)}), \hat{y}^{(n)})}{\partial w_i} = \frac{\partial E_{SV}(h(x^{(n)}), \hat{y}^{(n)})}{\partial h(x^{(n)})} \cdot \frac{\partial h(x^{(n)})}{\partial w_i}$$

$$\frac{\partial E_{SV}(h(x^{(n)}), \hat{y}^{(n)})}{\partial h(x^{(n)})} = \begin{cases} 0 & 1 - \hat{y}^{(n)} h(x^{(n)}) \leq 0 \\ -\hat{y}^{(n)} & 1 - \hat{y}^{(n)} h(x^{(n)}) > 0 \end{cases}$$

$$\frac{\partial h(x^{(n)})}{\partial w_i} = x_i^{(n)}$$

$$\frac{\partial \frac{1}{2} \|w\|^2}{\partial w_i} = w_i$$

(2). 代码实现:

- 初始化: 数据中标签的取值为{0,1}, 而为了训练模型需要把其改为{-1,1}
由观察得特征的取值范围较大, 则参数应该初始化为较小的数

```

1  def __init__(self, filename_train, filename_test, alpha, time):
2
3      # Parameters:
4      # filename_train (str): 训练数据的csv文件名字。
5      # filename_test (str): 测试数据的csv文件名字。
6      # 实例变量:
7      # - self.X_train_std: 标准化后的训练数据特征值。
8      # - self.Y_train: 训练数据的标签。{-1,1}
9      # - self.X_test_std: 标准化后的测试数据特征值。
10     # - self.Y_test: 测试数据的标签。 {-1,1}
11     # - self.alpha: 训练学习率
12     # - self.time: 训练迭代数
13     # - self.w_arg: 模型参数值
14     # - self.Loss_list: 存储每次训练周期的loss

```

```

15 # Note:
16 # - 在X的第一维度为添加的偏置因子
17 # - Y更改为{-1, 1}
18
19 X_train,Y_train=loadCSVfile1(filename_train)
20 X_test,Y_test=loadCSVfile1(filename_test)
21 sc=StandardScaler()
22 sc.fit(X_train)
23 X_train_std=sc.transform(X_train)
24 X_test_std=sc.transform(X_test)
25
26 new_column = np.ones((X_train_std.shape[0], 1))
27 self.X_train_std=np.hstack((new_column, X_train_std))
28
29 new_column_2= np.ones((X_test_std.shape[0], 1))
30 self.X_test_std=np.hstack((new_column_2, X_test_std))
31
32 self.Y_train=Y_train
33 self.Y_train[self.Y_train==0]==-1
34 self.Y_test=Y_test
35 self.Y_test[self.Y_test==0]==-1
36 self.time=time
37 self.alpha=alpha
38 self.w_num=self.X_train_std.shape[1]
39 self.w_arg=0.005*np.ones((self.w_num,1))
40 self.Loss_list=[]

```

- $L(w, b)$ 计算:

$$\text{公式: } L(w, b) = \sum_{n=1}^N E_{SV}(y^{(n)}h^{(n)}) + \frac{1}{2}||w||^2$$

```

1 def Loss(self):
2
3 # 算出所有训练样本的总loss，没有平均所以比较大
4 # return:
5 # - Loss: self.w_arg对应的当前模型对训练样本的总loss值
6
7 Y_train_Mat=np.mat(self.Y_train).transpose()
8 tm=1-
np.multiply(Y_train_Mat,np.dot(self.X_train_std,self.w_arg)) #计算 1-
yh(x)
9 w_tm = sum(x**2 for x in self.w_arg) / 2 #计算正则因子
10 Loss=0
11 Loss+=w_tm
12 print(tm.shape)
13 for i in range(0,tm.shape[0]):
14     if tm[i][0]>0:
15         Loss+=(tm[i][0])[0,0] #计算ESV(yh(x))的和
16     return Loss

```

- 单次梯度下降: 只要用矩阵向量运算

$$grad = - \sum_{n=1}^N y^{(n)} x_i^{(n)} + w_i \text{ 且 } 1 - y^{(n)} h(x^{(n)}) > 0$$


```

1  def Gradient_Descent(self):
2
3  # 对参数self.w_arg每个元素梯度下降更新一次
4  # Notes:
5  # - grad: 计算加入了正则因子的梯度, 没有平均
6
7  Y_train_Mat=np.mat(self.Y_train).transpose()
8  tm=1-
np.multiply(Y_train_Mat,np.dot(self.X_train_std,self.w_arg)) #计算 1-
yh(x)
9  tm[tm<=0]=0 #当1-yh(x)<=0时没有惩罚, 置为0
10 tm[tm>0]=1
11 tm=-np.multiply(tm,Y_train_Mat) #每个1-yh(x)>0的惩罚为 -yx_i
12 grad=np.dot(self.X_train_std.transpose(),tm)+self.w_arg #计算梯度
13 self.w_arg=self.w_arg-self.alpha*grad #梯度下降

```

- 多次迭代训练模型, 记录每次迭代的loss, 且在测试集上测试模型的正确率:

```

1  def Iterate(self):
2
3  # 1. 梯度更新self.time次, 训练模型并且记录loss
4  # 2. 计算训练完的模型对测试数据正确率, 并且打印输出总测试个数和总正确数
5  # 1:
6  for i in range(0,self.time):
7      loss=self.Loss()
8      print("hinge_loss_train loss: ",loss)
9      self.Loss_list.append(loss)
10     self.Gradient_Descent()
11 # 2:
12 predict=np.dot(self.X_test_std,self.w_arg)
13 correct_sum=0
14 for i in range(0,len(self.Y_test)):
15     if(predict[i][0]>=0):
16         if self.Y_test[i]==1:
17             correct_sum+=1
18     elif(predict[i][0]<0):
19         if self.Y_test[i]==-1:
20             correct_sum+=1
21     print("total test:",self.X_test_std.shape[0])
22     print("after hinge_loss_train,the correct_sum in test:
",correct_sum)

```

- 实验结果: 训练出的模型在测试集上的正确率较高,训练过程中loss不断下降

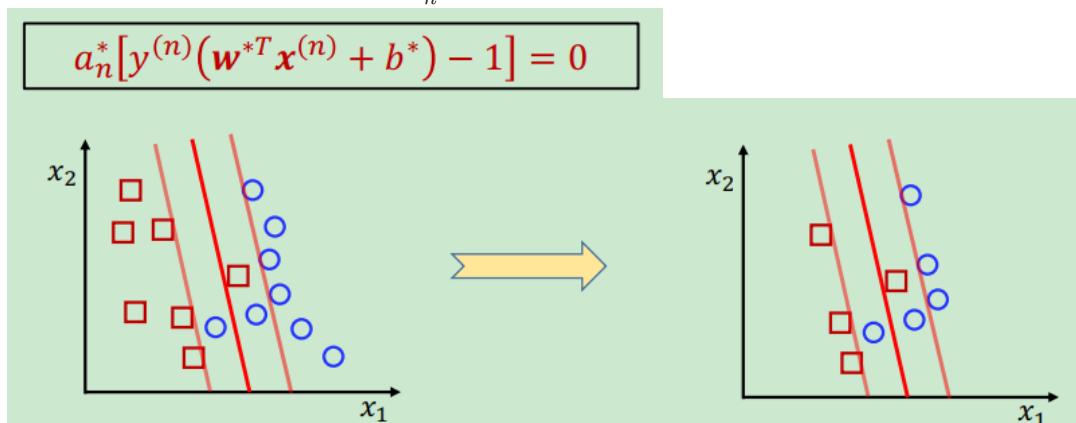

```

hinge loss训练时长: 4.9577 秒
total test: 2115
after hinge_loss_train,the correct_sum in test: 2113
hinge_loss_train loss: [[221271.16353214]]
(12665, 1)
hinge_loss_train loss: [[198894.6509174]]
(12665, 1)
hinge_loss_train loss: [[178827.77812678]]
(12665, 1)
hinge_loss_train loss: [[160807.81157938]]
(12665, 1)
hinge_loss_train loss: [[144597.62945343]]
(12665, 1)
hinge_loss_train loss: [[130006.57459131]]
(12665, 1)
hinge_loss_train loss: [[116856.71150755]]
(12665, 1)

```

(3). 采用 hinge loss 的线性分类模型和 SVM 模型之间的关系:

- hinge loss 和 SVM: svm (软间隔最大余量非线性分类器) 在非过拟合的情况下可能无法用一个超平面就将样本完全分类, 需要一个 ξ 给予每个样本可以偏移超平面的松弛量。 ξ 就是当样本可以被超平面正确分类时为0, 而不能被超平面正确分类时, 为了达到最小松弛量, ξ 取 $1 - y^{(n)}h(x^{(n)})$, 这和hinge loss的定义相似
- 采用hinge loss的线性分类器和SVM模型: 采用hinge loss的线性分类模型为了梯度下降将会将所有训练样本数据纳入计算范围
而SVM模型由于最优化问题: 计算 a_n^* 时只计算考虑在间隔范围内的样本点



3: 采用 hinge loss 线性分类模型和 cross-entropy loss 线性分类模型比较:

(1). cross-entropy loss交叉熵线性模型理论 (逻辑回归):

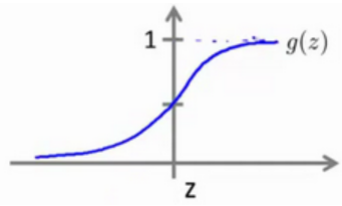
- 逻辑回归的假设函数公式: θ^T 为要训练的参数, x 为训练数据数据输入的特征

h_θ 可以近似看为得到的分类为1概率, 当它大于等于0.5时判别分类标签为1, 小于0.5时判别分类标签为0

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$g(z)$ 的图像如下: 由此易得 $\theta^T x \geq 0, h_\theta(x) \geq 0.5$,即分类为1



- $\theta^T(x)$ 的选取: 有两个特征值分别为 $x_1(\text{age})$ 和 $x_2(\text{薪水})$ 还有一个隐藏的 $x_0 = 1$ 作为偏置因子
此次实验实现了两个模型, 根据测试模型2的效果更好, 以下皆为模型2的分析

1. $\theta^T(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

$$2. \quad \theta^T(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2^2$$

- 参数 θ 的训练：使用梯度下降法， $J(\theta)$ 为代价函数， α 为学习率，公式为：

$$\theta_i = \theta_i - \alpha \frac{\delta J(\theta_i)}{\delta \theta_i}$$

注意：每次同时对所以特征参数迭代

梯度下降法求能使代价函数 min 的参数

梯度下降法: $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

代入

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \left[-\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(1 + e^{-\theta^T x^{(i)}}) - (1 - y^{(i)}) \log(1 + e^{\theta^T x^{(i)}})] \right] \\ &= -\frac{1}{m} \sum_{i=1}^m \frac{x_j^{(i)} e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} y^{(i)} - (1 - y^{(i)}) \frac{x_j^{(i)} e^{\theta^T x^{(i)}}}{1 + e^{\theta^T x^{(i)}}} \\ &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{x_j^{(i)}}{e^{\theta^T x^{(i)}} + 1} - (1 - y^{(i)}) \frac{x_j^{(i)}}{1 + e^{\theta^T x^{(i)}}} \\ &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} x_j^{(i)} - (1 - y^{(i)}) x_j^{(i)} e^{\theta^T x^{(i)}}}{e^{\theta^T x^{(i)}} + 1} \\ &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} (1 + e^{\theta^T x^{(i)}}) - e^{\theta^T x^{(i)}}}{1 + e^{\theta^T x^{(i)}}} x_j^{(i)} \\ &= -\frac{1}{m} \sum_{i=1}^m [(y^{(i)} - \sigma(\theta^T x^{(i)})) x_j^{(i)}] \\ &= \frac{1}{m} \sum_{i=1}^m [\sigma(\theta^T x^{(i)}) - y^{(i)}) x_j^{(i)}] \end{aligned}$$

- 代价函数 $J(\theta)$ 和损失函数 $Cost(h_{\theta}(x^i), y)$ 的定义:

$$Cost(h_\theta(x^i), y) = -y * \log(h_\theta(x)) - (1 - y) * \log(1 - h_\theta(x))$$

$$J(\theta) = \sum_{i=1}^m Cost(h_{\theta}(x^i), y^i)$$

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

m : 数据的数量 $x^{(i)}$: 第 i 个数据的特征 $x_j^{(i)}$ 第 i 个数据的第 j 个特征 $y^{(i)}$: 第 i 个数据的标签值

- 加入正则化后的参数训练：防止过拟合，为每个参数加上代价惩罚（一般不对偏置因子对应参数 θ_0 ）做惩罚，

$$J(\theta) = \sum_{i=1}^m (Cost(h_{\theta}(x^{(i)}) + \frac{\lambda}{m} \theta_j)$$

$$\theta_j = \theta_j - \alpha [\sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)} + \frac{\lambda}{m} \theta_j]$$

(2). 代码实现:

- 初始化: 注意样本的标签取值为{0,1}和hinge loss改为{-1,1}不同, 这里由于 $h(x) = wx + b$ 最后要经过sigmoid函数, 则 w 的初始值可以不用取得太小, 可以全都取为1

```

1  def __init__(self,filename_train,filename_test,alpha,time):
2
3      # Parameters:
4      # filename_train (str): 训练数据的csv文件名字。
5      # filename_test (str): 测试数据的csv文件名字。
6      # 实例变量:
7      # - self.X_train_std: 标准化后的训练数据特征值。
8      # - self.Y_train: 训练数据的标签。{0,1}
9      # - self.X_test_std: 标准化后的测试数据特征值。
10     # - self.Y_test: 测试数据的标签。 {0,1}
11     # - self.alpha: 训练学习率
12     # - self.time: 训练迭代数
13     # - self.w_arg: 模型参数值
14     # - self.Loss_list: 存储每次训练周期的loss
15     # Note:
16     # - 在X的第一维度为添加的偏置因子
17
18     X_train,Y_train=loadCSVfile1(filename_train)
19     X_test,Y_test=loadCSVfile1(filename_test)
20     sc=StandardScaler()
21     sc.fit(X_train)
22     X_train_std=sc.transform(X_train) #特征值标准化
23     X_test_std=sc.transform(X_test)  #特征值标准化
24
25     new_column = np.ones((X_train_std.shape[0], 1)) #加上偏置因子
26     self.X_train_std=np.hstack((new_column, X_train_std))
27
28     new_column_2= np.ones((X_test_std.shape[0], 1)) #加上偏置因子
29     self.X_test_std=np.hstack((new_column_2, X_test_std))
30
31     self.Y_train=Y_train
32     print(self.Y_train.shape)
33     self.Y_test=Y_test
34     self.time=time
35     self.alpha=alpha
36     self.w_num=self.X_train_std.shape[1]
37     self.w_arg=np.ones((self.w_num,1))
38     self.Loss_list=[]

```

- Loss计算:

```

1  def Loss(self):
2
3  # 算出所有训练样本的总loss，没有平均所以比较大
4  # return:
5  # - loss[0]: self.w_arg对应的当前模型对训练样本的总loss值
6
7  y=self.Y_train
8
9  Hypothesis_predict=self.sigmoid(np.dot(self.X_train_std,self.w_arg))
10     epsilon = 1e-5
11     loss=-(np.dot(y,np.log(Hypothesis_predict+epsilon))+np.dot((1-
12 y),np.log(1-Hypothesis_predict+epsilon)))
13     w_tm = sum(x**2 for x in self.w_arg) / 2
14     loss[0]+=w_tm
15     return loss[0]

```

- 梯度下降:

$$grad = (Y - sigmoid(X * W))^T * X^T$$

```

1  def Gradient_Descent(self):
2
3  # 对参数self.w_arg每个元素梯度下降更新一次
4  # Notes:
5  # - grad: 计算加入了正则因子的梯度，没有平均
6
7  Y_train_Mat=np.mat(self.Y_train).transpose()
8  grad=-np.dot((Y_train_Mat-
9 self.sigmoid(np.dot(self.X_train_std,self.w_arg))).transpose(),self.X_train_std).transpose()
10     grad+=self.w_arg
11     self.w_arg=self.w_arg-self.alpha*grad

```

- 训练模型并且在测试集上测试:

```

1  def Iterate(self):
2
3  # 1. 梯度更新self.time次，训练模型并且记录loss
4  # 2. 计算训练完的模型对测试数据的正确率，并且打印输出总测试个数和总正确数
5
6  # 1:
7  for i in range(0,self.time):
8      #阶段性计算损失
9      loss=self.Loss()
10     print("cross_entropy_train loss: ",loss)
11     self.Loss_list.append(loss)
12     self.Gradient_Descent()
13
14 #2:
15 predict=np.dot(self.X_test_std,self.w_arg)
16 correct_sum=0
17 for i in range(0,len(self.Y_test)):
18     if(predict[i][0]>=0 and self.Y_test[i]==1):
19         correct_sum+=1
20     elif(predict[i][0]<0 and self.Y_test[i]==0):
21         correct_sum+=1

```

```
21         print("total test:",self.x_test_std.shape[0])
22         print("after cross_entropy_train,the correct_sum in test:
",correct_sum)
```

- 结果：训练时间较长，但是在测试集上准确率高, 训练过程中loss不断下降

```
cross_entropy_loss训练时长： 12.4765 秒
total test: 2115
after cross_entropy_train,the correct_sum in test:  2113
cross_entropy_train loss:  216451.11259586038
cross_entropy_train loss:  194813.46540430858
cross_entropy_train loss:  175341.30651276527
cross_entropy_train loss:  157826.71484693963
cross_entropy_train loss:  142075.98404453497
cross_entropy_train loss:  127913.97247323432
cross_entropy_train loss:  115160.02263454418
cross_entropy_train loss:  103674.06341771333
```

(3). 和 hinge loss 线性分类模型比较:

- 训练时间上：cross-entropy loss 线性分类模型比hinge loss长，cross-entropy loss模型计算中涉及log exp等复杂计算，而hinge loss只有简单的加减乘除
- 训练效果：
 - 在该测试样本上一样的超参数设置下正确率相等且都很高，这侧面体现了数据可以被线性分类。
 - hinge loss线性模型鼓励模型在正确分类的同时，最大化决策边界的间隔，而cross entropy loss线性模型则不考虑优化决策边界的间隔
 - hinge loss线性模型在计算梯度和loss时会看重分类错误和距离决策边界太近的样本点，其他样本点的惩罚都为0，而cross entropy loss关注计算所有样本点，所有训练样本点都有loss

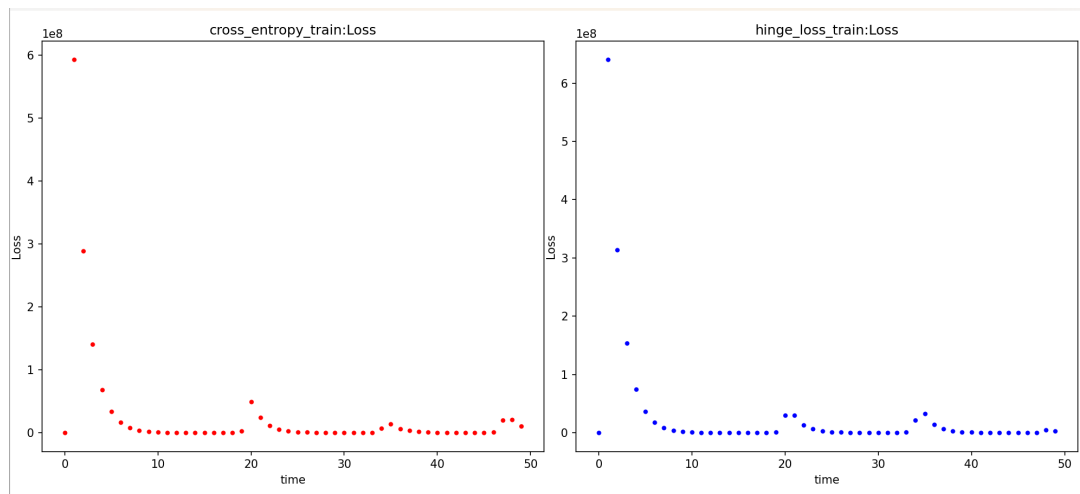
三：实验结果

1. 不同超参数设置：

- 学习率：由于在样本的loss计算和梯度计算时没有对训练样本平均，所以梯度和loss都会比较大，则学习率不宜设的过大，不然会引起震荡

- 设为0.3时，引起震荡：可以看到loss-time图中，无论是hinge loss还是cross entropy loss都随着迭代次数增加有明显震荡

```
cross_entropy_train loss: 99498.93919411008
cross_entropy_train loss: 48088.80434779373
cross_entropy_train loss: 23248.933166441806
cross_entropy_train loss: 11456.315676255299
cross_entropy_train loss: 6526.9761438539135
cross_entropy_train loss: 25164.154453452127
cross_entropy_train loss: 2940536.229923964
cross_entropy_train loss: 49292917.27188314
cross_entropy_train loss: 24103463.119005077
hinge_loss_train loss: [[111725.21329681]]
(12665, 1)
hinge_loss_train loss: [[54228.33808916]]
(12665, 1)
hinge_loss_train loss: [[26388.00429854]]
(12665, 1)
hinge_loss_train loss: [[13234.21598292]]
(12665, 1)
hinge_loss_train loss: [[6747.90240255]]
(12665, 1)
hinge_loss_train loss: [[6897.46080591]]
(12665, 1)
hinge_loss_train loss: [[875113.64072795]]
(12665, 1)
hinge_loss_train loss: [[29723810.73586074]]
(12665, 1)
hinge_loss_train loss: [[29903868.87149192]]
```

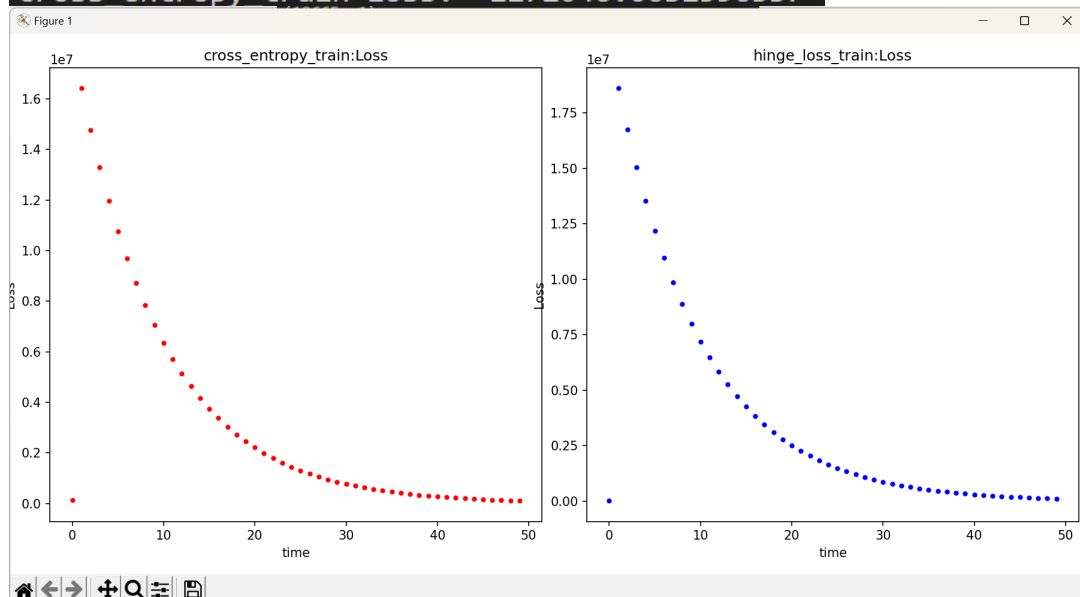


- 设为0.05时: loss随着迭代次数的增加一直在下降

```

hinge_loss_train loss: [[273862.14053153]]
(12665, 1)
hinge_loss_train loss: [[246195.02753827]]
(12665, 1)
hinge_loss_train loss: [[221271.16353214]]
(12665, 1)
hinge_loss_train loss: [[198894.6509174]]
(12665, 1)
hinge_loss_train loss: [[178827.77812678]]
(12665, 1)
hinge_loss_train loss: [[160807.81157938]]
(12665, 1)
hinge_loss_train loss: [[144597.62945343]]
(12665, 1)
hinge_loss_train loss: [[130006.57459131]]
(12665, 1)
hinge_loss_train loss: [[116856.71150755]]
(12665, 1)
cross_entropy_train loss: 5142917.846482186
cross_entropy_train loss: 4627626.665942866
cross_entropy_train loss: 4163869.809798137
cross_entropy_train loss: 3746537.9340782706
cross_entropy_train loss: 3370995.932594235
cross_entropy_train loss: 3033132.171538703
cross_entropy_train loss: 2729130.506767851
cross_entropy_train loss: 2455596.0638074465
cross_entropy_train loss: 2209441.774574745
cross_entropy_train loss: 1987937.0082637875
cross_entropy_train loss: 1788615.1331291283
cross_entropy_train loss: 1609252.4899025213
cross_entropy_train loss: 1447875.0948981256
cross_entropy_train loss: 1302680.9903059856
cross_entropy_train loss: 1172046.0631596337

```



- 迭代次数: 一开始设为迭代100次时, 发现后期loss下降的曲线越来越平缓, 在50次左右趋向稳定了, 改为迭代50次, 效果和迭代100次差不多, 但是减少计算和训练时间

2. 所有实验结果:

样本适合在线性平面内分类, 线性核svm和基于hinge loss或者cross-entropy loss的线性分类模型训练正确率都较高

```
高斯核SVM训练时长: 2.7910 秒
gauss Misclassified samples: 8
gauss Accuracy: 0.996217
线性核SVM训练时长: 0.5421 秒
linear Misclassified samples: 1
linear Accuracy: 0.999527
cross_entropy_loss训练时长: 17.2984 秒
total test: 2115
after cross_entropy_train,the correct_sum in test: 2113
hinge_loss_train loss: [[116856.71150755]]
hinge_loss训练时长: 3.1423 秒
total test: 2115
after hinge_loss_train,the correct_sum in test: 2113
```