```
中山大学移动信息工程学院本科生实验报告
                    (2017年秋季学期)
课程名称:人工智能
                  专业方向
                                     学号
   年级
                                                  姓名
   1501
                移动 (互联网)
                                    15352005
                                                  蔡景韬
一、实验题目
 • 逻辑回归模型 Logistic Regression Model
二、实验内容
1. 算法原理
 • 软分类、硬分类
   o 软分类
      ■ 基于概率模型的分类模型,输出不同类对应的概率,最后取概率最高的类别作为分类结果
      ■ 概率模型:形式为P(X|y),即在学习过程中,y未知,训练后模型得到的输出是x的一系列值的概率
      ■ 如,NB分类、逻辑回归分类
   o 硬分类
      ■ 基于非概率模型的分类模型,输出决策函数的决策结果
      ■ 非概率模型:形式为决策函数,即输入x到输出y的一个映射,且输出唯一
      ■ 如,决策树、PLA、KNN算法等
 • 逻辑回归(从PLA到LR)
   o 通过训练得到权重矩阵,用以表示各列属性与结果的相关程度
      ■ 权重为正,说明该属性与"是"的结果正相关,且值越大,正相关程度越大,正影响越大
      ■ 权重为负,说明该属性与"是"的结果负相关,且绝对值越大,负相关程度越大,负影响越大
   ο 基本的预测公式
      ■ 将权重矩阵与属性值内积(相应的对属性值进行加权),并与某个阈值作差
      ■ 上述表示可写出表达式: y = (\sum_{i=1}^d W_i X_i) - threshold
   o 公式的数学处理
      ■ 为了简化表达式,我们可以把阈值当作W_0,且对X矩阵多添加X_0=1(阈值可以当作一个属性且每
       一个训练文本共用一个阈值,于是可以把阈值属性值都取1)
      • 表达式可以重写为: y = \sum_{i=0}^d W_i X_i = W^T X
```

o PLA算法\_预测公式 ■ 对于上文得到的公式, PLA使用sign()函数对公式结果进行处理, 从而得到PLA的预测公式:  $f(x) = sign(W^T X)$ ■ 当预测公式得到正结果,则预测为正结果 ■ 否则预测为负结果 o PLA算法\_更新公式 对于权值的更新, 使用梯度下降法进行更新 (之前写过详细推导, 这里就不赘述) ■ 使用向量模型理解有, X矩阵是一个高维的向量, 权重矩阵是另一个高维向量, X矩阵与权重矩阵的 内积的正负可以表征两个向量之间的夹角 ■ 于是可以得到更新公式:  $W^{p+1} = W^p + y_i X_i$ o LR算法\_预测公式 ■ 对于上文得到的基本预测公式, 预测的y值范围为 $[-\infty, +\infty]$ ■ 如果我们想得到概率,明显需要对得到的y值进行处理 ■ 考虑logit函数 ■  $log(\frac{p}{1-p}) = y$ , 其中p是正结果的概率 ■ 验证其合理性: ■ 当p趋近于1时(即正结果概率近乎为1),可以看到,y趋近于+∞,根据PLA算法,知道 此时是预测为正结果的 ■ 当p趋近于0时(即正结果概率近乎为0),可以看到,y趋近于+∞,根据PLA算法,知道 此时是预测为负结果的 ■ 根据上文的logit函数,于是可以求解正结果概率p  $log(rac{p}{1-p}) = W^T X \quad \Rightarrow \quad p = rac{1}{1+e^{-W^T X}}$ ■ 可以得到LR算法的概率公式:  $f(x) = P(y|x) = p^y(1-p)^{1-y}$ , 其中y表示x对应的分类标签 • y = 1, f(x) = P(1|x) = p $\blacksquare \quad \exists \ y = 0, \ \ f(x) = P(0|x)1 - p$ o LR算法\_更新公式 对于权值,使用极大似然法进行更新(PLA使用的是最小二乘法) ■ 似然函数:  $egin{aligned} likelihood &= \prod_{i=1}^M P(label|x_i) \ &= \prod_{i=1}^M p^y (1-p)^{1-y} \end{aligned}$ ■ 左右取-log有  $Err(W) = -log(likelihood) = -logigg(\prod_{i=1}^{M} p^y (1-p)^{1-y}igg)$  $=-\sum_{i=1}^{M}\left[y_{i}log(p)+(1-y_{i})log(1-p)
ight]$ 其中, $p=rac{1}{1+e^{-W^TX}}$ ■ 使用极大平均似然法进行更新,即量 Err(W)取最小时,似然函数达到最大值 ■ 根据高数知识,当*Err(W)*梯度为0时,函数取到最小值(关于连续可导性、二阶可微、凸函数的证 下列式子中,W为权重矩阵, $X_n$ 为第n行的数据集矩阵, $y_n$ 为第n行的label:  $abla Err(W) = rac{\partial Err(W)}{\partial W} = -\sum_{n=1}^{N} \left[ y_n rac{rac{\partial p}{\partial W}}{p} + (1-y_n) rac{-rac{\partial p}{\partial W}}{1-p} 
ight]$  $egin{aligned} &\Rightarrow hinspace 1 + e^{-W^T X_n} = u, \; \mathbb{N} \; p = rac{1}{u} \ &= -\sum_{n=1}^N \left[ y_n rac{1}{p} rac{-rac{\partial u}{\partial W}}{u^2} + (1-y_n) rac{-1}{1-p} rac{-rac{\partial u}{\partial W}}{u^2} 
ight] \end{aligned}$  $= -\sum_{n=1}^{N} \left[ y_n \frac{1}{p} \frac{-1}{u^2} (e^v) \frac{\partial v}{W} + (1 - y_n) \frac{-1}{1 - p} \frac{-1}{u^2} (e^v) \frac{\partial v}{W} \right]$  $=-\sum_{n=1}^{\infty}\left[y_nrac{1}{p}rac{-1}{u^2}(e^v)(-X_n)+(1-y_n)rac{-1}{1-p}rac{-1}{u^2}(e^v)(-X_n)
ight]$  $=-\sum_{i=1}^{N}\left[y_{n}\,rac{1}{p}(p)(1-p)(X_{n})+(1-y_{n})rac{-1}{1-p}(p)(1-p)(X_{n})
ight]$ 

 $=-\sum_{n=1}^{N}\left[y_{n}(1-p)(X_{n})+(1-y_{n})(-p)(X_{n})
ight]$  $=\sum_{n=1}^N(p-y_n)X_n$ ■ 又 $\nabla Err(W)$ 是关于矩阵的N元函数,难以求得 $\nabla Err(W) = 0$  ,所以使用数值分析中的迭代方法求解 零点,则可以得出权值的更新公式:  $W_{t+1} = W_t - rac{\eta}{M} 
abla Err(W_t)$  $=W_t-rac{\eta}{M}\sum_{i=1}^N(p-y_n)X_n$  $V_t = W_t - rac{\eta}{M} \sum_{n=1}^N \left[ \left(rac{1}{1+e^{-W^TX_n}} - y_n
ight) X_n 
ight]$ ■ 注意事项 ■ 如果e的指数太大或太小,如 $e^{18}$ 或 $e^{-18}$ ,算出来的 $\frac{1}{1+e^{-W^TX_n}}$ 要么是 $\frac{1}{\infty}=0$ ,或者是 $\frac{1}{1+0}=1$ ■ 要控制点乘运算 $W^TX_n$ 的大小 ■ 初始化W不能太大 ■ 学习率n不能太大 o LR算法\_正则化 ■ 正则化的概念 ■ 当模型参数很多,而我们可用的数据非常少时,极易出现过拟合的问题。为此,就引入了正则 ■ 其目的是使得参数空间受到一定的限制。

"Just right"

 $J(w)\cong J(w)+\lambda \|w\|\left\{egin{aligned} L2 & o rac{1}{2}\lambda \|w\|_2^{-2} \ L1 & o \lambda \|w\|_1 \end{aligned}
ight. 
onumber \left. \|x\|_p := \left(\sum_{i=1}^n |x_i|^p
ight)^{rac{1}{p}}$ 

■ 而L1就比较突兀,可能会直接使得某些参数的取值为0.。

 $Err(W) = -rac{1}{M}\sum_{i=1}^{M}\left[y_ilog(p) + (1-y_i)log(1-p)
ight] + rac{1}{2}\lambda\sum_{j=1}W_j^2$ 

■ 在误差函数 Err(W)中引入正则项, 重写误差函数有

■ 总结起来就是: L1会引入稀疏性, 而L2会充分利用更多的特征。

■ 如上图中, 维数过多、数据量过少导致了过拟合现象, 于是在成本函数中引进正则项(惩罚

■ 正则项的加入改变了参数的学习规则,它在参数的每一步更新中都新增了一个系数参量。

■ L2对参数的正则更加的平滑,即它限制了参数空间,但对参数的影响是平滑的

High bias

(underfit)

■ 正则化一般分为L1正则化与L2正则化

项)限制其参数空间

■ L1正则和L2正则的差别

■ 本次实验使用L2正则化

o LR与DT/RF(决策树/随机森林)的比较:

乏探查局部结构的内在机制。

■ 但是LR可以在线更新且提供有用的概率信息。

■ 两种方法都很快且可扩展。 ■ 在正确率方面, RF比LR更优。

■ 两者的算法差别

■ 总结:

o 计算某一列属性的文本梯度和

for i:文本数Rows

根据梯度公式求解每一行文本的梯度,并求和:

void getTrainVal(string str="IS NOT TEST"){

void doTrain(int iterations, double eta=1, double lambda=0){

int Size = Data.size(),TrainSize=Size\*0.75 ;

Train.assign(Data.begin(),Data.end()) ;

Train.assign(Data.begin(),Data.begin()+TrainSize) ; Val.assign(Data.begin()+TrainSize,Data.end()) ;

if ( str=="IS NOT TEST" ){

( iterations-- ){ eta \*= 0.995 ; update(eta,lambda);

void update(double eta, double lambda){ vector<double> W\_new(W.size()); int TrainSize = Train.size();
for ( int i=0 ; i<Col ; i++ )</pre>

o 使用0-3代码优化

执行1000次的情况下

■ 没有使用优化的情况下

■ 使用优化的情况下

o 用二维数组代替vector

使用正则化与动态学习率

1. 实验结果展示示例(使用小数据集)

 $=2-\eta \times 0.0554$ 

 $=1-\eta \times 0.0758$ 

• 程序执行结果 (可以看到是正确的答案)

2. 评测指标展示即分析(如果实验题目有特殊要求,否则使用准确率)

W矩阵初始化为1,执行1000次,学习率η=0.1,正则参数λ=0.1,动态学习率每次乘以0.995

可以得到准确率曲线(准确率-选代次数)(下降的曲线是使用动态学习率时的学习率数值)

800

迭代次数

正则

0.7785

1000

1200

普通

0.779

正则+动态学习

动态学习率

正则

普通

D:\GT2\As student\AI

ost time: Oms

0.75 0.7 0.65 0.6 0.55 0.5 0.45

0.4

0.35 0.3

0.25 0.2

0.15 0.1 0.05

正则+动态

0. 7815

四、思考题

• 批梯度下降

• 随机梯度下降

o 优点:

o 缺点:

200

• 从图像可以看到, 收敛速度: 正则 > 正则+动态学习 > 普通

• 从收敛准确率表格可以看到,收敛准确率:正则+动态学习 > 普通 > 正则

1. 如果把梯度为0作为算法停止的条件,可能存在怎样的弊端?

3. 思考批梯度下降和随机梯度下降这两种优化方法的优缺点

o 批梯度下降算法在更新权值矩阵的时候要遍历整个数据集

1) 训练过程可能十分漫长;

o 随机梯度下降的思想是根据每个单独的训练样本来更新权值

1) 由于每次仅仅采用一个样本来迭代, 训练速度很快

2) 由于每次迭代方向不同,可能会跳出局部最优,而收敛到全局最优

随机梯度下降法用于仅仅用一个样本决定梯度方向,导致解很有可能不是最优。

的, 因为梯度值会越来越小, 它和固定的学习率相乘后的积也会越来越小

原因: 动态学习率到后面可以一步一步的逼近最优解, 而正则能够避免过拟合

原因: 动态学习率到中期学习率较小, 所以收敛速度也慢

可以得到最终的收敛准确率表格

5.01167 99446 992419

=0.9924

• 得到的W矩阵

 $w_2^{new} = w_2^{old} - \eta \sum \left[ \left( \frac{e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}}{1 + e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}} - y \right) x_2 \right]$ 

=1.9945

三、实验结果及分析

• 准确率优化

训练集

步长: η = 0.1

cost time: 393191ms

time: 56974ms

第j列属性

label

3

3. 关键代码截图

LR算法

• 划分训练集与验证集

■ 进行L2正则化

• 动态学习率

Size  $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ 

High variance (overfit)

■ 重新推导误差函数的梯度 下式中的w是W矩阵的一个维度  $abla Er ilde{r}(w) = 
abla \left[ Err(w) + rac{1}{2} \lambda w^2 
ight]$  $= \nabla Err(w) + \lambda w$ ■ 则Err(W)的梯度重写为 $\nabla Err(W) = \nabla Err(W) + \lambda W$ ■ 更新公式重写为  $W_{t+1} = W_t - \eta igg(rac{1}{N} \sum_{n=1}^N \left[ig(rac{1}{1 + e^{-W^T X_n}} - y_nig) X_n
ight] + \lambda W_tigg)$  $=(1-\eta\lambda)W_t-rac{\eta}{N}\sum_{n=1}^N\left[ig(rac{1}{1+e^{-W^TX_n}}-y_nig)X_n
ight]$ • 与上文相同,需要保证 $W^TX_n$ 足够小,所以 $\lambda$ 也要较小 o 在迭代过程前半部分学习率应该设置较大, 因为初始化的W是随机给的, 可能离最优解很远 o 在迭代过程的后半部分学习率应该设置较小, 因为已经接近最优解附近, 此时应该减小学习率(步 长), 让其不至于错过最优解, 可以一步一步的靠近最优解 ο 算法实现: 每次迭代都将 η 乘以0.998 • 逻辑回归的优缺点 o LR有很多方法来对模型正则化,如L1正则化、L2正则化,从而避免了欠拟合和过拟合的问题。 o LR有很好的概率解释,且很容易利用新的训练数据来更新模型。 o 比起NB的条件独立性假设, LR不需要考虑样本是否是相关的。

■ 逻辑回归对数据整体结构的分析优于决策树,而决策树对局部结构的分析优于逻辑回归。

■ 逻辑回归对极值比较敏感,容易受极端值的影响,而\*决策树对极值有很好的抗干扰性。

线性联系在实践中有很多优点:简洁,易理解,可以在一定程度上防止对数据的过度拟合。

■ 决策树由于采用分割的方法, 所以能够深入数据细部, 但同时就失去了对全局的把握。

■ 同时由于切分,样本数量不断萎缩,所以无法支持对多变量的同时检验。

■ 一个分层一旦形成,它和别的层面或节点的关系就被切断了,以后的挖掘只能在局部中

■ 逻辑回归,始终着眼整个数据的拟合,所以对全局把握较好。但无法兼顾局部数据,或者说缺

■ 逻辑回归擅长辨识线性关系,而决策树对线性关系的把握较差。

2. 伪代码 • 划分训练集与验证集 训练集 = 数据集的前75% 验证集 = 数据集的后25% LR算法 o 训练W矩阵 while iterations --动态学习率: eta=eta\*0.995 更新W矩阵: update(eta,lambda) o 更新W矩阵 for i:属性列数Cols 根据矩阵更新公式更新W矩阵,lambda是正则项系数: 3 W\_new[i] = (1-eta\*lambda)\*W[i] - 1/TrainSize\*eta\*文本梯度和(i)

sum += ( 1/(1+exp(- W矩阵·第i个文本矩阵) ) - 第i个文本第j列属性值 ) \* 第i个文本的

W.swap(W\_new); double gradSum(int col){ double ans=0,TrainSize=Train.size(); for ( int i=0 ; i<TrainSize ; i++ ){
 double dot\_product = dotProduct(Train[i]) ;</pre> ans += ( (double)1/(double)(1+exp(-dot\_product)) - Train[i][Col] ) \* Train[i][col] ; return ans ; double dotProduct(const vector<double>& data){ double sum=0; for ( int j=0 ; j<Col ; j++ )</pre> sum += data[j]\*W[j] ; n sum ; 4. 创新点&优化 • 时间优化

D:\GT2\As student\AI\_Lab\Lab5\15352005\_caijingtao.exe

D:\GT2\As student\AI Lab\Lab5\15352005 caijingtao.exe

执行1000次迭代时能减少2s,但后来加了正则化等优化之后就被vector赶上来了,这里就不贴了

W\_new[i] = (1-eta\*lambda)\*W[i] - eta\*gradSum(i)/TrainSize ;

```
• 迭代一次计算
              \tilde{\mathbf{w}}^{\mathsf{T}}\tilde{\mathbf{x}}_{A} = \begin{pmatrix} -5 \\ 2 \\ 1 \end{pmatrix}^{\mathsf{T}} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = w_{0} \times 1 + w_{1} \times x_{A}^{1} + w_{2} \times x_{A}^{2} = -5 \times 1 + 2 \times 0 + 1 \times 1 = -4
               \tilde{\mathbf{w}}^{\mathsf{T}}\tilde{\mathbf{x}}_{\scriptscriptstyle R} = -5 \times 1 + 2 \times 1 + 1 \times 1 = -2
               \tilde{\mathbf{w}}^{\mathrm{T}}\tilde{\mathbf{x}}_{C} = -5 \times 1 + 2 \times 3 + 1 \times 3 = 4
               \tilde{\mathbf{w}}^{\mathrm{T}}\tilde{\mathbf{x}}_{n} = -5 \times 1 + 2 \times 4 + 1 \times 3 = 6
             w_0^{new} = w_0^{old} - \eta \sum \left[ \left( \frac{e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}}{1 + e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}} - y \right) x_0 \right]
                             = -5 - \eta \left[ \left( \frac{e^{-4}}{1 + e^{-4}} - 0 \right) \times 1 + \left( \frac{e^{-2}}{1 + e^{-2}} - 0 \right) \times 1 + \left( \frac{e^{4}}{1 + e^{4}} - 1 \right) \times 1 + \left( \frac{e^{6}}{1 + e^{6}} - 1 \right) \times 1 \right]
                              =-5-\eta \times 0.1167
                              =-5.0117
             w_1^{new} = w_1^{old} - \eta \sum \left| \left( \frac{e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}}{1 + e^{\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}} - y \right) x_1 \right|
                             =2-\eta \left[ \left( \frac{e^{-4}}{1+e^{-4}} - 0 \right) \times 0 + \left( \frac{e^{-2}}{1+e^{-2}} - 0 \right) \times 1 + \left( \frac{e^{4}}{1+e^{4}} - 1 \right) \times 3 + \left( \frac{e^{6}}{1+e^{6}} - 1 \right) \times 4 \right]
```

 $=1-\eta \left[ \left( \frac{e^{-4}}{1+e^{-4}} - 0 \right) \times 1 + \left( \frac{e^{-2}}{1+e^{-2}} - 0 \right) \times 1 + \left( \frac{e^{4}}{1+e^{4}} - 1 \right) \times 3 + \left( \frac{e^{6}}{1+e^{6}} - 1 \right) \times 3 \right]$ 

• 迭代很难得到精确值,如果能够得到满足精度的近似值,迭代就可以停止 • 如果把算法停止的条件设置为梯度为0,则 o 算法可能无法停止, 因为数据集的梯度可能永远无法收敛到0 o 数据集的梯度可以收敛到0,但所用的时间太长,没有意义 2.η 的大小会怎么影响梯度下降的结果? 给出具体的解释, 可视化的解释最好, 比如图形 展示等 如果 η 值太大,则每次迭代就有可能出现超调的现象,会在极值点两侧不断发散,最终损失函数的值是越 变越大, 而不是越来越小, 且还可能出现该曲线不断震荡的情形。 o 如下图图示,由于学习率太大,导致错过了最优解,且由于下次的逼近又会跳回左侧,所以会出现不断 振荡的情形 如果 η 值太小,这该曲线下降得很慢,甚至在很多次迭代处曲线值保持不变,最终导致在规定的迭代次数 中梯度无法下降到一个较好的结果

1.0	1,000
0.5	***************************************
0.0	
-0.5	P2 P3
-1.0	
-1.5	The state of the s
-2.0	2.0 -1.5 -1.0 -0.5 0.0 0.5 1.0 1.5 2.0
	缺点:
	在训练数据特别庞大的时候,可能会出现
	在州外数据有为此人的时候,「形云山地

2) 如果误差曲面上有多个局极小值, 那么不能保证这个过程会找到全局最小值。

o 批梯度下降每次移动的方向都是确定的,一直更新后在合适的学习率情况下可以收敛到一个局部最小点

对于收敛速度来说,由于随机梯度下降法一次迭代一个样本,导致迭代方向变化很大,不能很快的收敛 到最优解。