

《人工智能实验》

实验报告

(期中项目)

学 院 名 称 : 数据科学与计算机学院

专业 (班级) : 16 计科 2 班

学 生 姓 名 : 李新锐

学 号 : 15323032

时 间 : 2018 年 10 月 21 日

一、数据预处理

数据清理

向量化

二、KNN

创新点

实验结果

三、朴素贝叶斯

算法原理

伪代码

实验创新

代码展示

实验结果

四、回归树

算法原理

数据划分

终止条件

伪代码

代码

算法改进

实验结果

五、随机森林

Bagging

伪代码

从Bagging到随机森林

伪代码

代码

实验结果

六、逻辑斯特回归

算法原理

伪代码

代码

实验结果

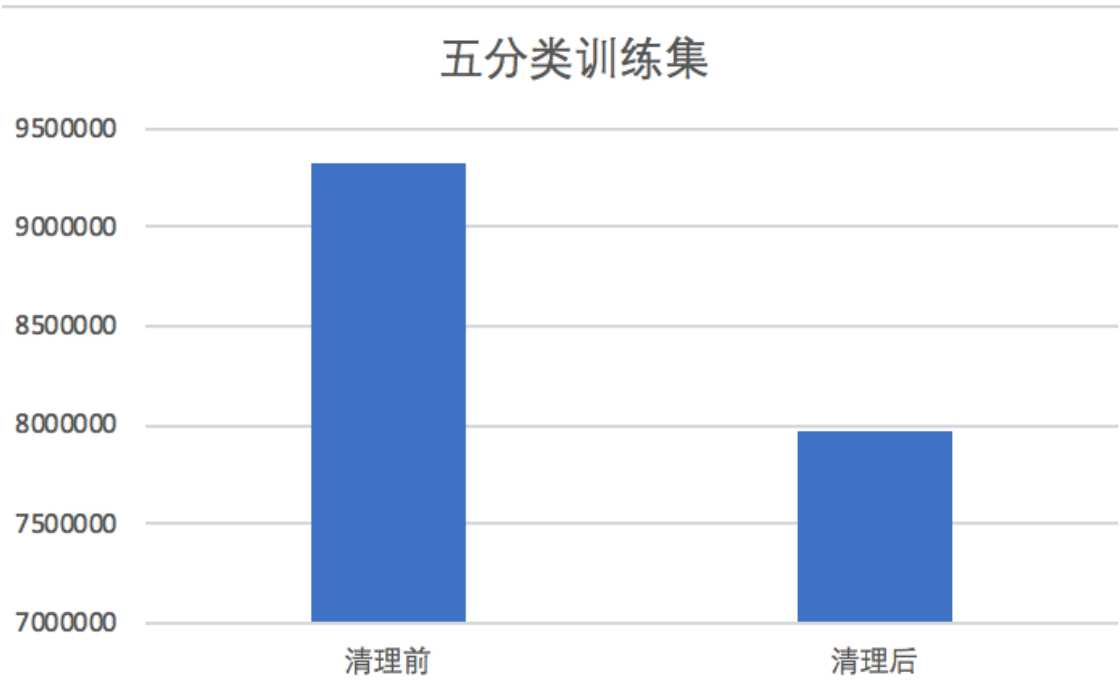
一、数据预处理

数据清理

对于二分类和五分类的训练样本，由于它们是原始的自然语言文本，因此在向量化之前，我对它们进行了较为全面的数据清洗工具，包括了以下几个步骤：

1. 使用`nlk.tokenize`进行分词
2. 使用了`isalpha`函数去除非单词内容。

该步骤去除了文本中的数字、标点、特殊符号、HTML标签等非英文内容。以五分类数据集为例，清理前后训练集单词总量从93万下降到79万，降幅达15。



3. 进行拼写检查

自然语言文本中有许多拼写错误，检测并更正这些错误能够减少文本中的噪声。

我首选尝试的是PyEnchant工具和aspell-en字典，然而人工检查发现它把很多人名识别为了错误。因此改用了autocorrect包，人工检查发现效果不错。

以五分类数据集为例，拼写检查结果如下：

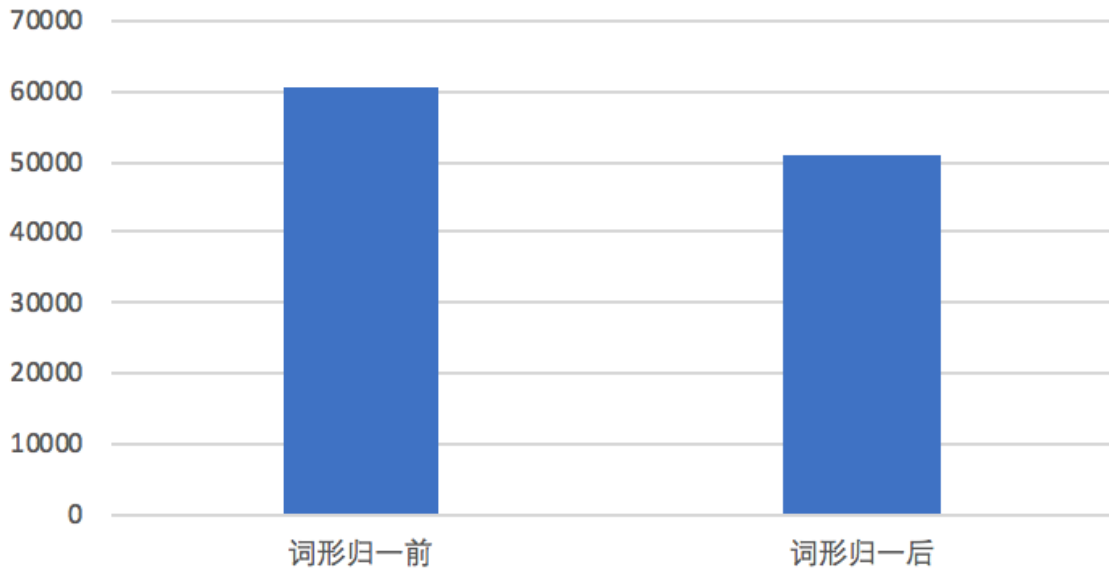
更正数目	更正率
187672	2.345%

4. 词形归一化（包括大写转小写）

英语文本的一个特点是同一单词可能有着不同的形式。如take taken taking。在自然语言处理中这些词形信息往往不重要，通过词形归一化能够降低不同的单词数目，起到降维的效果。我使用的归一化工具是nltk.stem.WordNetLemmatizer。

在五分类训练集上词形归一化前后文本中不同的单词数目从60760下降到了51169。相当于降维了近1万维。

五分类训练集



向量化

本次实验中我主要使用了Onehot矩阵和doc2vec矩阵。其中前者是在实验一中实现的，后者调用了gensim库。该方法和实验课上所讲的word2vec类似，区别在于训练集中每个句子被分配一个id，在训练模型的输入层中和句子中所有单词放在一起训练，相当于在训练的过程中用到了整个句子的信息。训练的结果被当做整个句子的向量表示。

最终得到的向量表示如下：

数据集	维度
二分类-onehot	62751
二分类-doc2vec	50
五分类-onehot	51170
五分类-doc2vec	50

二、KNN

创新点

本次实验中，我首先尝试了直接使用lab1中实现的KNN算法，由于样本个数过大，预测速度非常缓慢。因此作出一处改进：将计算L2范数的过程改为矩阵运算。原理如下：

设 X 为 $N \times M$ 的训练集， T 为 $P \times M$ 的测试集。

目标是得到 $P \times N$ 的距离矩阵 D 。 D 中元素 $D_{i,j}$ 是测试集向量 T_i 和训练集中向量 X_i 的距离。则有

$$D_{i,j} = \sqrt{\sum_k^m (X_{i,k} - T_{j,k})^2}$$

代码实现：

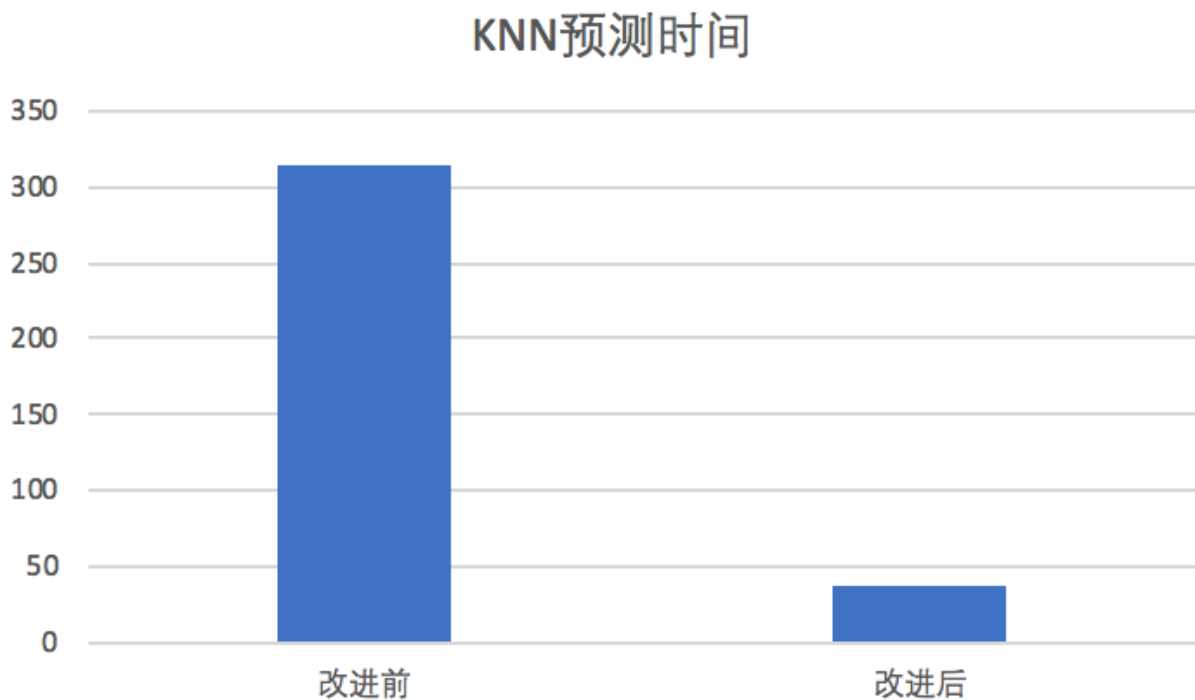
```

1  #trainSet, testSet, dists分别对应原理说明中的X,T,D矩阵
2  #trainSum: N x 1
3  trainSum = np.sum(np.square(trainSet), axis=1)
4  #testSum: P x 1
5  testSum = np.sum(np.square(testSet), axis=1)
6  #t0:P x N
7  t0 = np.dot(testVec, trainSet[0].T)
8  #dist: P x N
9  dists = np.sqrt(-2 * t0 + testSum.reshape(-1, 1) + trainSum)

```

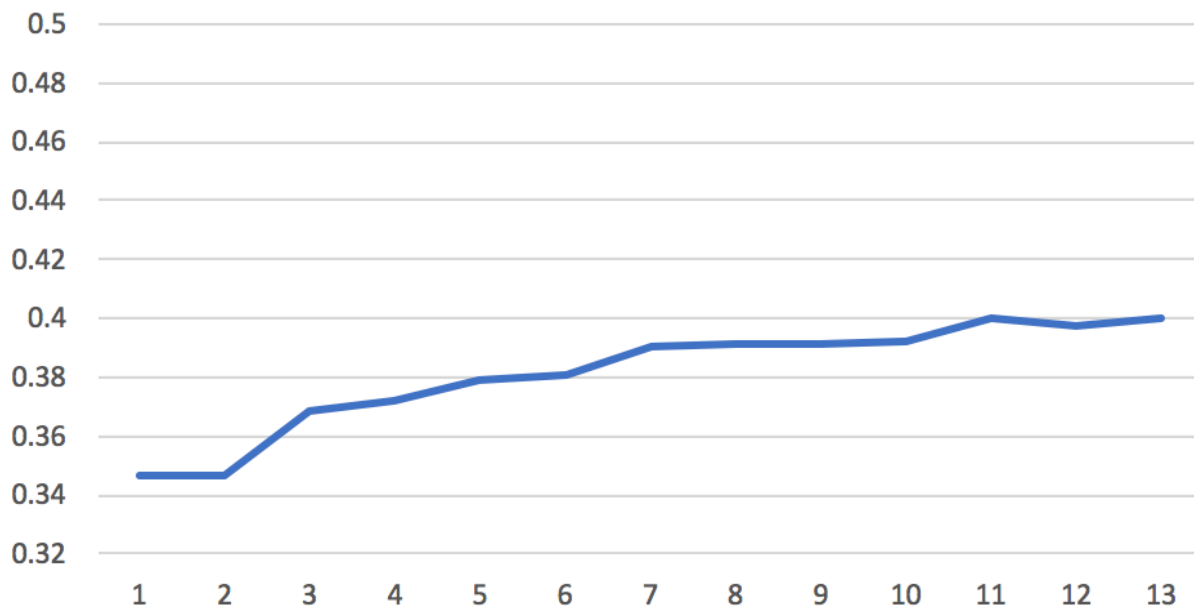
实验结果

我使用了doc2vec矩阵，在二分类训练集中，选取后4000条作为验证集。采用上述的改进方式前后，预测时间从5分降低到34秒左右。



但缺点是，上述改进限制了只能使用二范数。因此可调参数只有K值。在K值从1到13的过程中，二分类 分类准确度变化如下：

K值-二分类准确率



可见，准确度均为达到50%。可能原因是doc2vec的降维并不是一个保范数的过程，50维的doc2vec样本空间中的样本直接的空间相对远近和原始空间中各样本之间的相对远近是不一致的。

因此，我尝试改用onehot矩阵，但二分类的onehot矩阵有62751维，用KNN无法在合理的时间内处理。因此放弃了并寻求更好的方法。

三、朴素贝叶斯

算法原理

本次实验中我尝试的第二种算法是朴素贝叶斯算法。

朴素贝叶斯是一种分类算法，基于的是概率论中贝叶斯公式。

假设我们有 N 个样本，每个样本有 m 个特征。样本分为 K 各类别 C_1, C_2, \dots, C_K 。分类问题要求的是在输入为 X 的条件下输出为 Y_k 的概率。由条件概率公式得：

$$P(Y_k|X) = \frac{P(Y_k X)}{P(X)}$$

而 $P(Y_k X)$ 可以写作 $P(X|Y_k)P(Y_k)$ ，同时由全概率公式知

$$P(X) = \sum_k P(Y_k)P(X|Y_k)$$

因此我们就得到了贝叶斯公式：

$$P(Y_k|X) = \frac{P(X|Y_k)P(Y_k)}{\sum_k P(X|Y_k)P(Y_k)}$$

而如何利用该公式进行分类呢，方法如下：

1. 计算样本属于各个类别的先验分布

$$P(Y = C_K)$$

2. 计算条件概率 $P(X|Y = C_K)$ 。假设 X 的各个特征之间是独立的。可得

$$P(X|Y = C_K) = P(X_1 = x_1|Y = C_k)P(X_2 = x_2|Y = C_k) \\ \dots * P(X_m = x_m|Y = C_k)$$

3. 由于贝叶斯公式的分母 $P(X)$ 是一样的，因此预测分类结果 C^* 是使得分母最大化的 Y 的取值。因此贝叶斯分类的问题就变为了求解：

$$\begin{aligned} & \operatorname{argmax}_{C_k} P(X|Y = C_K) * P(Y = C_K) \\ &= \operatorname{argmax}_{C_k} P(X_1 = x_1|Y = C_k)P(X_2 = x_2|Y = C_k) \\ & \quad \dots * P(X_m = x_m|Y = C_k) * P(Y = C_K) \end{aligned}$$

4. 最后要解决的问题如何计算有关 X 的各个特征 m 的条件概率 $P(X_m|Y = C_k)$ 。对于特征取值离散的oneHot矩阵，可以直接统计频率构造伯努利分布。然而使用doc2vec进行向量化后，各个维度 X_s 都是连续值。因此处理方法是：假设各个特征都服从正态分布，通过样本计算出均值和方差，即可拟合出这个分布。

伪代码

oneHot版本：

```
1 # 记训练集输入为D，标签集为Y，要预测的输入为X
2 # 计算K个先验概率
3 for k in 1...K:
4     P_Y[k] = count(Y == k) / count(D)
5 # 按照类别划分输入集
6 for k in 1...K:
7     D_parted[k] = filter(Y == k, D)
8 # 统计每一维度特征在K种类别上的伯努利分布
9 PXY = [[1-sum(attr)/N, sum(attr)/N] for attr in zip(*D_parted[k])]
10
11 # 计算输入X为K种类别的概率
12 P = P_Y
13 for k in 1...K:
14     for i in 1...m:
15         P[k] *= PXY[k][i][X[i]] #P(X_i = X_i | y_k)
16 # 取概率最大的一种作为输出
17 return max_argument(P)
```

doc2Vec版本与oneHot版本的区别在于

```
1 # 第八行改为
2 gauss = ([Guass(mean(attr), var(attr)) for attr in zip(*D_parted[k])])
3 # 统计每一维度特征在K种类别上的均值和方差，并依此构造高斯分布
4
5 # 第15行改为
6 P[k] *= gauss[k][i](X[i])
7 # 计算高斯分布在X[i]点的概率
```

实验创新

由于onehot矩阵的维度非常高，在代码第15行连续计算乘法会导致最终得到的概率被舍为0.因此改为取对数求加法。

```
1 P[k] += log(PXY[k][i][X[i]])
```

然而 $PXY[k][i][X[i]]$ 的值可能为0，会导致算出nan的情况。因此又要对PXY的计算进行改进，将第9行改为

```
1 1-(1+ sum(attr))/(N+2), (1+sum(attr))/(N+2)
```

来避免这种情况。

代码展示

```
1  ## D_parted[k] 对应分类为k的样本集
2  D_parted = list(range(K))
3  for k in range(2):
4      ## 通过filter函数筛选
5      D_parted[k] = np.array(list(filter(lambda v: v[-1] == k, D)))
6      ## 计算先验概率
7      P_Y[k] = len(D_parted[k]) / len(D)
8
9  ## 定义高斯分布类
10 class Guass:
11     # 均值
12     mean = 0
13     #标准差
14     stdev = 0
15     def __init__(self, stats):
16         self.mean, self.stdev = stats
17     # 计算x处的概率
18     def P(self, x):
19         t1 = np.exp(-(np.power(x-self.mean,2)/(2*np.power(self.stdev,2))))
20         return (1 / (np.sqrt(2*np.pi) * self.stdev)) * t1
21
22
23 guass = list(range(K))
24 # 统计每一维度特征在K种类别上的均值和方差，并依此构造高斯分布
25 for k in range(K):
26     guass[k] = [Guass((np.mean(attr), np.std(attr))) for attr in zip(*D_parted[k])]
27     # 删除最后一列(Y值)构造的高斯分布
28     del guass[k][-1]
29
30 # 统计每一维度特征在K种类别上的伯努利分布，和上面的代码一样，为了使用tqdm库显示运行进度条才拆开了写，所有看起来长一点
31 onehotPXY = list(range(K))
32 for k in range(K):
33     l = list(range(len(X[0]) + 1))
```



```

34     z = zip(*D_parted[k])
35     for i, attr in enumerate(tqdm_notebook(z)):
36         l[i] = (1-(sum(attr)+1)/(len(D_parted[k])+2), (sum(attr)+1)/(len(D_parted[k])))
37     del l[-1]
38     onehotPXY[k] = l
39
40 # 预测部分
41 def predict(X):
42     # 使用deepcopy防止修改P_Y
43     _PY = deepcopy(P_Y)
44     # 计算输入X为K种类别的概率
45     for k in range(K):
46         for i in range(M):
47             ret = np.log(guass[k][i].P(X[i]))
48             _PY[k] += ret
49     #取概率最大的一种作为输出
50     return np.argmax(_PY)

```

实验结果

采用K交叉验证

使用doc2vec仅有31%的正确率。推测原因是：输入向量各维度是正态分布的假设和doc2vec向量的实际不同。

使用onehot时正确率非常好，5-交叉和6-交叉验证的正确率均为90.88%。

K值	正确率
5	90.8875%
6	90.8875%

提交rank后，在测试集上的准确率为85.38，比在验证集上低5%左右。这说明验证集上的数据与测试集分布式有一定差别的。由于

```
16337237_0.txt Ok 85.38333333333334
```

四、回归树

由于朴素贝叶斯算法缺少调参空间，在得到测试结果后，我接着尝试了改进实验二中实现的决策树。当时我只实现了处理离散数据的分类树，但doc2vec矩阵各个维度上都是连续值，因此我实现了回归树算法。算法原理如下：

算法原理

数据划分

将测试数据依据某个特征 a ，以 s 为切分点进行切分。求解

$$\operatorname{argmin}_{a,s} \left[\sum_{x_i < s} (\operatorname{avg}_1 - y_i)^2 + \sum_{x_i \geq s} (\operatorname{avg}_2 - y_i)^2 \right]$$

其中 avg_1 是 s 左边数据的平均值

avg_2 是 s 及右边数据的平均值

每个区域的输出即是平均值

终止条件

1. 误差函数的值差距小于阈值

$$\sum_{i \in D} (\operatorname{avg} - y_i)^2 < \operatorname{tolerance}$$

2. 选择最好的分割方法,不纯度减少量仍低于阈值

$$\sum_{x_i < s} (\operatorname{avg}_1 - y_i)^2 + \sum_{x_i \geq s} (\operatorname{avg}_2 - y_i)^2 - \sum_{i \in D} (\operatorname{avg} - y_i)^2 < \operatorname{tolerance}$$

3. 数据量小于阈值

$$|D| < \operatorname{tolerance}$$

伪代码

```
1  # 依据训练样本集D和特征集A训练回归树
2  # errFunc可以采用方差
3  train(D, errFunc):
4      # 训练集D的X, Y分开
5      (X, Y) = D
6      # 记avg为当前节点数据集的结果的平均数
7      avg = sum(Y) / |D|
8      '''
9      结束条件有3种, 分别是:
10         1. 误差函数的值差距小于阈值
11         2. 选择最好的分割方法, 不纯度减少量仍低于阈值
12         3. 数据量小于阈值
13     '''
14     if errFunc(D) < tolerance1 or |D| < tolerance2
15         # 设置当前节点为叶节点, 结果为mode
16         node.isLeaf = True
17         node.Y = avg
18         return node
19
20     bestA = -1, bestErr = errFunc(D)
21     # 遍历一种特征
22     for a in A:
```

```

23     # 对每一种该特征的取值
24     for s in X[a]:
25         # 划分数据集为两份
26         D_l, D_r = split_data(s)
27         # 计算新误差
28         newErr = errFunc(D_l) + errFunc(D_r)
29         #记录最佳划分
30         if deltaErr < bestErr:
31             bestErr = deltaErr
32     # 选择最好的分割方法,不纯度减少量仍低于阈值
33     if bestA == -1:
34         node.isLeaf = True
35         node.Y = avg
36         return node
37     #递归训练子节点
38     node.child.append(train(D_l))
39     node.child.append(train(D_r))

```

代码

```

1  // 实际代码中考虑了深度限制, 因此多了int depth (当前深度), int level (最大深度) 两个参数
2  void RegressionTree::train_worker(RegressionTree::Node *node, const ErrFunc_t &errFunc,
3  int depth, int level)
4  {
5      //训练集D的X, Y分开
6      auto* pA = &node->A;
7      auto* pD = &node->D;
8      auto D_Y = pD->row(-1);
9      // 记avg为当前节点数据集的结果的平均数
10     auto avg = D_Y.sum() / D_Y.size();
11     //当前误差
12     auto err = errFunc(D_Y, avg);
13     //cout << "In level " << depth << endl;
14     /*
15     结束条件有3种, 分别是:
16     1. 误差函数的值差距小于阈值
17     2. 选择最好的分割方法,不纯度减少量仍低于阈值
18     3. 数据量小于阈值
19     */
20     if(err < RegressionArgs::err_tolerance || depth > level)
21     {
22         node->isLeaf = true;
23         node->Y = avg;
24         return;
25     }
26
27     int bestA = -1;
28     double bestS = INFINITY;
29     double bestErr = err;

```

```

29 // 遍历一种特征
30 for(const auto& a : range(0, pA->size()))
31 {
32     auto vals = pD->row(a);
33     std::set<double> vals_set;
34     auto n = vals.size();
35     //实际代码中, 只抽样该特征的部分取值
36     for(Eigen::Index i = 0; i < 10; ++i)
37     {
38         vals_set.insert(vals[RandLib::uniform_rand(0, n - 1)]);
39     }
40     for(auto s: vals_set)
41     {
42         //划分数据集
43         auto[D_l, D_r] = split_data(node->D, a, s);
44         auto Y_l = D_l.row(-1);
45         auto Y_r = D_r.row(-1);
46         auto avg_l = Y_l.sum() / Y_l.size();
47         auto avg_r = Y_r.sum() / Y_r.size();
48         auto err_l = errFunc(Y_l, avg_l);
49         auto err_r = errFunc(Y_r, avg_r);
50         //计算误差变化
51         auto new_err = err_l + err_r;
52         if(new_err < bestErr)
53         {
54             bestA = a;
55             bestS = s;
56             bestErr = new_err;
57         }
58     }
59 }
60 //cout << "Best cut: new_err = " << bestErr << " , A = " << bestA << " , S = " <<
bestS << endl;
61 //选择最好的分割方法,不纯度减少量仍低于阈值, 则结束递归
62 if(bestA == -1)
63 {
64     node->isLeaf = true;
65     node->Y = avg;
66     return;
67 }
68 auto[D_l, D_r] = split_data(node->D, bestA, bestS);
69 //判断左右子节点是否有数据集为空
70 if(D_l.empty())
71 {
72     node->isLeaf = true;
73     node->Y = avg;
74     return;
75 }
76 if(D_r.empty())

```

```

77     {
78         node->isLeaf = true;
79         node->Y = avg;
80         return;
81     }
82     //否则递归训练子节点
83     node->C = bestA;
84     node->S = bestS;
85     node->ch_l = new Node(D_l, *pA);
86     node->ch_r = new Node(D_r, *pA);
87     train_worker(node->ch_l, errFunc, depth+1, level);
88     train_worker(node->ch_r, errFunc, depth+1, level);
89 }

```

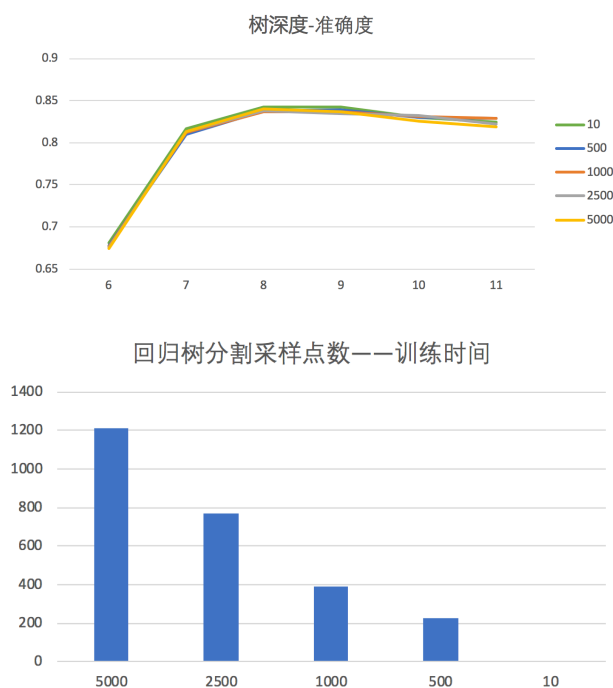
算法改进

如上述代码所示，在实际实现回归树时，并不是遍历某一特征所有取值，而只是随机抽取部分点，从而降低训练决策树的时间。这种分割可能无法找到最优分割点，如下实验结果所示：即便仅抽取10个分割点，在树达到一定深度时，分类效果不比抽取更多分割点差。

实验结果

- 测试方法：K=5的K交叉验证
- 测试对象：二分类问题（测试集大小：19200，验证集大小：4800）
- 测试输入：50维的doc2vec矩阵
- 图像中不同直线含义：在选取特征分割点时分别采样10，500，1000，2500，5000个点
- 横轴含义：树的最大深度
- 纵轴含义：准确度

可见，在树深度为8至9时有着最佳的分类准确度。其中使用5000个分割采样点的直线在深度为8时有准确度最大值84.0362，然而其他方法也与他相差无几。但带来的好处是近似线性的训练速度提升。



(注：10对应的训练时间为4，所有时间单位为秒)

五、随机森林

在实现了回归树算法后，我不断调整树的最大深度和采样点树，然而似乎训练效果不能继续提升了。因此我尝试采用了随机森林算法。

算法原理如下：

Bagging

随机森林算法的基础是Bagging。该算法对大小为 N 的原数据集进行 N 次放回的采样。由于有放回这一特点，导致有些点会被采样多次，有些点无法采样到，因此能够减少算法的过拟合。

伪代码

```
1  rets = 0
2  for i in range(M):
3      # D是样本集
4      sub_D = sample_D(D)
5      # A是特征集
6      tree = CART.train(sub_D, A)
7      ret = tree.vaild(vaild_set)
8      rets += ret
9  rets /= M
```

从Bagging到随机森林

由于回归树是一个贪婪算法，即便使用Bagging，得到的树在结构上仍然相似，结果高度相关。为了减少模型间的相关度，随机森林算法做出了进一步改进，在特征集上也进行抽样。

伪代码

```
1  rets = 0
2  for i in range(M):
3      # D是样本集
4      sub_D = sample_D(D)
5      # A是特征集
6      sub_A = sample_A(A)
7      tree = CART.train(sub_D, sub_A)
8      ret = tree.vaild(vaild_set)
9      rets += ret
10 rets /= M
```

代码

```
1  // 对样本集trainSet采样
2  matrix_view<double> BaggingRegressTree::sample_D(const matrix_view<double>& trainSet)
```

```

3 {
4     //v存储被采样的样本的下标
5     std::vector<Eigen::Index> v;
6     //使用集合类型自动去除重复
7     std::set<Eigen::Index> s;
8     auto n = trainSet.cols();
9     //随机采样n次。RandLib::uniform_rand是我封装的均匀分布函数
10    for(Eigen::Index i = 0; i < n; ++i)
11        s.insert(RandLib::uniform_rand(0, n - 1));
12    //返回采样到的样本
13    for(auto i : s)
14        v.push_back(i);
15    return matrix_view<double>(trainSet, v);
16 }
17 // 从大小为n的特征空间中抽取k个特征
18 Vec<Eigen::Index> BaggingRegressTree::sample_A(size_t n, size_t k)
19 {
20     Vec<Eigen::Index> ret;
21     //若k大于等于n直接返回
22     if(k >= n)
23         return range(Eigen::Index(0), Eigen::Index(n));
24     std::set<Eigen::Index> sel;
25     //循环抽取直到抽取到k个不同的特征
26     while(sel.size() < k)
27     {
28         int s = RandLib::uniform_rand(0, n - 1);
29         sel.insert(s);
30     }
31     //返回被抽取的特征
32     for(auto i : sel)
33         ret.push_back(i);
34     std::sort(ret.begin(), ret.end());
35     return ret;
36 }
37 //训练随机森林，M代表森林中树的数量，k代表抽取的特征数
38 void BaggingRegressTree::train(const ErrFunc_t& errFunc, int M, int k, int maxLevel)
39 {
40     //建立M棵树
41     for(int i = 0; i < M; ++i)
42     {
43         //样本集采样
44         matrix_view<double> sub_D = sample_D(trainSet);
45         //特征集采样
46         Vec<Eigen::Index> sub_A_idx = sample_A(featureCount, k);
47         //储存到随机森林类中
48         sub_A_idx.push_back(-1);
49         forests_A.push_back(sub_A_idx);
50         sub_D.select_row(sub_A_idx);
51         //建树

```

```

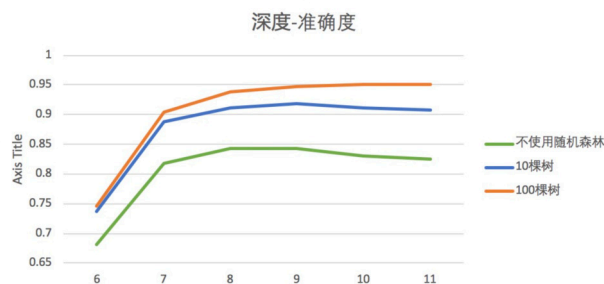
52     RegressionTree* tree = new RegressionTree(sub_D);
53     forests.push_back(tree);
54     //训练
55     tree->train(errFunc, maxLevel);
56 }
57 }

```

实验结果

- 测试方法：K=5的K交叉验证
- 测试对象：二分类问题（测试集大小：19200，验证集大小：4800）
- 测试输入：50维的doc2vec矩阵
- 图像中不同直线含义：不使用随机森林，随机森林建10棵树，建100棵树
- 横轴含义：树的最大深度
- 纵轴含义：准确度

如图可视：建100棵树的随机森林算法取得了出乎意料的好结果，在 $K \geq 9$ 验证集分类准确度稳定达到95%！



六、逻辑斯特回归

最后，我尝试并实现了逻辑斯特回归算法。

算法原理

逻辑斯特回归是一种分类算法，它的原理是在样本空间中划分线性边界 $a = W^T X$ ，再通过激活函数

$$h(x) = \frac{1}{1 + e^{-x}}$$

将边界两端的点映射到0, 1两种分类上。

从上面的描述可以看到，逻辑斯特回归模型的训练实际上就是去找到一个最佳的 W 矩阵。根据最大似然估计，可以定义模型的损失函数为

$$C(W) = - \sum_{i=1}^n (y_i \log h(W^T x_i) + (1 - y_i) \log (1 - W^T x_i))$$

通过梯度下降法可以对 W 的每个维度 j 分别进行更新

$$\begin{aligned}
 W_{new}^j &= W^j - \mu \frac{\partial C(W)}{\partial W^j} \\
 &= W^j - \mu \sum_{i=1}^n [(h(W^T X_i) - y_i) X_i^j] \\
 &= W^j - \mu X^T (h(XW) - Y)
 \end{aligned}$$

公式中 X_i 和 y_i 是大小为 n 的样本集中的样本

更新的结束条件可以设置为：

1. 限制最大迭代次数 max_epoch
2. W 的变化小于阈值
3. $C(W)$ 的变化小于阈值

伪代码

```

1  # W初始化为M+1维的全1向量,
2  W = ones(M+1)
3  C = Cost(W)
4  epoch = 0
5  while True:
6      # err[i] 是梯度下降法公式中h(W^T * X_i) - y_i这一项
7      err = h(X * W) - Y
8      # 更新W
9      W = W - u * X.T * err
10     #判断结束条件
11     epoch += 1
12     if epoch > max_epoch or delta(W) < min_change or delta(C) < min_change:
13         break

```

代码

```

1  # sigmoid函数
2  def h(x):
3      return 1 / (1 + np.exp(-x))
4  # 最大似然cost函数
5  def cost(W):
6      ret = 0
7      # 简单的公式翻译
8      for i in range(N):
9          t0 = np.dot(W.T, X[i])
10         t1 = Y[i] * np.log(h(t0))
11         t2 = (1-Y[i]) * np.log(h(1-t0))
12         ret -= (t1 + t2)
13     return ret
14
15 # 最大迭代次数
16 max_epoch = 10000
17 # 记录最新的W
18 LastW = [0]

```

```

19 # 记录训练过程中的loss变化
20 loss = []
21 def train(W):
22     epoch = 0
23     while True:
24         # X: N x M
25         # W: M * 1
26         # err[i] 是梯度下降法公式中 $h(W^T * X_i) - y_i$ 这一项
27         err = h(X @ W) - Y
28         # 更新W
29         #X.T: M x N
30         #err: N x 1
31         #X.T @ err : M x 1
32         #学习率
33         u = 0.00001
34         deltaW = - u * (X.T @ err)
35         # 更新并记录新的W
36         W = W + deltaW
37         LastW[0] = W
38         # 记录新cost
39         newCost = cost(W)
40         loss.append(newCost)
41         print(epoch, newCost)
42         #判断结束条件
43         epoch += 1
44         if epoch >= max_epoch:
45             break

```

实验结果

- 测试方法：K交叉验证
- 测试对象：二分类问题（测试集大小：19200，验证集大小：4800）
- 测试输入：50维的doc2vec矩阵

在经过多次调整学习率之后，我的逻辑斯特回归模型仍未能收敛，而是陷入了loss约为8520的局部极小点。

K值	正确率
5	78.208%
6	78.208%