

# 数据分析及实践Exp4 Report

## 1. 实验描述

### 1.1 实验背景:

**英雄联盟**（League of Legends, LOL）是美国游戏开发商Riot Games（2011年由腾讯收购）开发和发行的一款多人在线战斗竞技游戏。

在游戏中，玩家扮演一个“召唤师”角色，每个召唤师控制一个拥有独特技能的“英雄”，并与一组其他玩家或电脑控制的英雄战斗。游戏的目标是摧毁对方的防御塔和基地。

召唤者峡谷是英雄联盟中最受欢迎的地图，在这种地图类型中，两队五名玩家竞争摧毁一个被称为基地的敌人建筑，这个建筑由敌人的队伍和一些防御塔护卫。每个队伍都希望保卫自己的建筑，同时摧毁对方的建筑。这些主要包括：

- Towers（防御塔）：每支队伍总共有11座防御塔
- Inhibitor（水晶）：每条道有一个水晶
- Dragon（大龙）
- Rift Herald（峡谷先锋）
- Baron Nasho（纳什男爵）
- Nexus（基地）

英雄联盟最具争议的元素之一，就是其严重的滚雪球效应。许多职业选手接受赛后采访时都提到其输赢都因为“滚雪球”，我们研究各方面各指标的数据，来看这些因素的发展是否真的影响了比赛的成败。在这个实验中，我们分析了8万场英雄联盟的排名比赛，尝试得出有效的结论。

### 1.2 实验要求:

#### Part1:

- 基于Exp3，手动实现一种分类算法（例如，决策树、KNN或者朴素贝叶斯。并参考Exp3的特征工程，测试算法在LOL数据集上的预测性能，撰写实验报告。

- 代码实现只允许使用 `numpy`、`pandas` 库和 `python` 内置库，不允许使用现有的机器学习库。
- 预测任务与实验三一致，以准确率作为评价指标。自行在 LOL 数据集上划分训练集和验证集（4:1 比例、交叉验证），汇报算法在验证集上的性能。

## Part2:

- 基于 Exp3，预测测试集中每个样本的比赛持续时间 `gameDuration`。
- 将数据集中一部分样本的标签，作为训练集，而另一部分样本作为测试集。
- 可以利用开源工具包，也可以参考 Exp3 的数据分析与特征工程。
- 以均方误差 (Mean-square Error) 作为评价指标。

### 1.3 数据集说明:

- `team1_win` 代表 `team1` 是否取得胜利
- `team1_firstBlood` 代表 `team1` 是否取得一血，其他类似特征同理
- `player1_championId` 代表玩家 1 选择的英雄 ID
- `player1_kills` 代表玩家 1 的击杀数，其他类似特征同理
- `player1 ~ player5` 属于 `team1`，`player6 ~ player10` 属于 `team2`
- `gameDuration` 代表游戏时长

## 2. Part1

### 2.1 算法选择

使用 **K近邻算法** (K-Nearest Neighbor, KNN)，基本思想是一个样本与数据集中的  $k$  个样本的“距离”最近，如果这  $k$  个样本中的大多数属于某一个类别，则该样本也属于这个类别。KNN 算法的关键在于  $k$  值的选取和“距离”计算的方式。

此处选择的“距离”为特征之差，选取的特征有：一血、一龙、一塔、一男爵、一先锋、一水晶、同时一血一龙、同时一龙一先锋、同时一男爵一先锋、同时一塔一水晶、场均击杀差、场均死亡差、场均助攻差、经济差、KDA 差共 15 个特征。

最初选用的距离计算方式为：各特征差的绝对值，乘以基于 Exp3 得出的经验权重（或与 `win` 的相关系数），最后再对各加权特征差求和，作为该样本间的距离。后来发现由于数据量过大，该计算方式会极大增加计算时间，所以选择以放弃少量精度为代价，来提高计算速度。

改进后的距离计算方式为：最初将训练集每一个样本的特征加权求和作为该样本的坐标，将多维问题降为一维问题。计算距离只需要对加权特征和作差即可。理论上该举措会降低精度，实际实验结果表明精度的变化在0.2%左右，以此为代价来极大地提高计算速度是值得的。

## 2.2 优化设计

KNN算法虽然是最简单的分类算法之一，但其计算量往往是巨大的。在最初版本中，在1/5的测试集上预测的时间高达900秒，这是不可接受的。于是对算法的计算进行优化。

- 首先在距离计算方式上进行改进，将多维计算变为一维计算，只需要计算样本间的加权特征和之差。
- 预处理时，生成含有每个样本特征的 dataframe，其中包含已归一化的特征值，所有特征值位于0~1之间，其中将一血类的 bool 值转化为 int 值，并放缩100倍（即0-0.01）。同时将所有特征的加权和（即坐标）也存入其中。
- 完成测试任务时，将测试集的 dataframe 转化为 array 来操作。
- 运用 array 的 lexsort() 函数来进行距离升序排序，并选取前k个。
- 运用 dataframe 的 value\_counts() 函数进行统计，选取出现次数最多的标签。

## 2.4 核心代码

```
1  #预处理后的数据数为75201，进行分组
2  h = int(75201/5)
3  H = [0,h,2*h,3*h,4*h,75201]
4  #随机打乱
5  df=df.reindex(np.random.permutation(df.index))
6  accuracy = []
7  for i in range(5):
8      #划分测试集与训练集
9      test = df[H[i]:H[i+1]]
10     train = pd.concat([df[0:H[i]], df[H[i+1]:75201]], axis=0)
11     #生成特征表
12     train = gettrain(train)
13     train_f = getfeature(train,train)
14     test_f = getfeature(train,test)
15     #测试数据个数、k值选取
16     right = 0
17     n = test.shape[0]
18     k = 37
```

```

19     #转化为array
20     t0 = np.array(train_f[['sum', 'win']])
21     t1 = np.array(test_f['sum']).tolist()
22     t2 = np.array(test_f['win']).tolist()
23     #开始测试
24     time_start=time.time()
25     for j in range(n):
26         t3 = abs(t0 - [t1[j],0])    #距离
27         index=np.lexsort([t3[:,0]])
28         temp = pd.DataFrame(t0[index, :][0:k])[1]    #排序取前k
29         个
30         result = temp.value_counts().index[0]    #统计取第一个
31         if result == t2[j]: #准确率统计
32             right += 1
33         accuracy += [right/n]
34         time_end=time.time()
35         print('Fold',i+1,':','准确率为: %.3f'%(accuracy[i]),'耗时:
36         %.6f s'%(time_end-time_start))
37
38     sum = 0
39     for i in accuracy:
40         sum = sum + i
41     avg_acc = sum/5
42     print('平均准确率为: %.3f'%(avg_acc))

```

## 2.3 结果分析

```

Fold 1 : 准确率为: 0.977 耗时: 81.188232 s
Fold 2 : 准确率为: 0.977 耗时: 81.058267 s
Fold 3 : 准确率为: 0.975 耗时: 80.908459 s
Fold 4 : 准确率为: 0.978 耗时: 81.801882 s
Fold 5 : 准确率为: 0.979 耗时: 81.275309 s
平均准确率为: 0.977

```

平均准确率为97.7%，是一个可观的结果，但每一次验证的平均时常约为81.2秒，尽管从900多秒到81.2秒已经经过了很多优化，但对于后续调参来说还是很不方便。

## 2.4 算法改进：类重心KNN

在传统的KNN算法中, 为了找到这个含有k个训练集中的样本的最近邻, 需要计算该未知样本点和所有训练集样本之间的距离, 然后从最小距离开始计样本数, 一直计算到有K个样本数为止。

有一种简化的算法称为类重心法<sup>[1]</sup>, 即将训练集中每类样本的重心求出, 然后判别位置样本点与各类重心的距离, 未知样本点距哪一类重心距离最近, 位置样本就属于哪一类。

本Part共两类——team1胜、team2胜。于是在求的特征表的基础上, 计算两类样本的重心, 即类特征均值, 计算测试样本与两个重心的距离, 极大的提升了计算时间, 并且对精度的影响并不大。

```
1 accuracy = []
2 for i in range(5):
3     #划分测试集与训练集
4     test = df[H[i]:H[i+1]]
5     train = pd.concat([df[0:H[i]], df[H[i+1]:75201]], axis=0)
6     #生成特征表
7     train = gettrain(train)
8     train_f = getfeature(train, train)
9     test_f = getfeature(train, test)
10    #重心计算
11    train_f_1 = train_f[train_f['win']==1]
12    train_f_0 = train_f[train_f['win']==0]
13    center_1 = train_f_1.sort_values(by = ['sum'])
    [(train_f_1.shape[0]//2):(train_f_1.shape[0]//2+1)]
    ['sum'].values[0]
14    center_0 = train_f_0.sort_values(by = ['sum'])
    [(train_f_0.shape[0]//2):(train_f_0.shape[0]//2+1)]
    ['sum'].values[0]
15
16    right = 0
17    n = test.shape[0]
18    #转化为array
19    t1 = np.array(test_f['sum']).tolist()
20    t2 = np.array(test_f['win']).tolist()
21    #开始测试
22    time_start=time.time()
23    for j in range(n):
24        if abs(center_1-t1[j]) <= abs(center_0-t1[j]): #与重心
            距离
```

```

25         result = 1
26     else:
27         result = 0
28         if result == t2[j]:
29             right += 1
30     accuracy += [right/n]
31     time_end=time.time()
32     print('Fold',i+1,':','准确率为: %.3f'%(accuracy[i]),'耗时:
%.6f s'%(time_end-time_start))
33
34     sum = 0
35     for i in accuracy:
36         sum = sum + i
37     avg_acc = sum/5
38     print('平均准确率为: %.3f'%(avg_acc))

```

结果:

```

Fold 1 : 准确率为: 0.978 耗时: 0.017894 s
Fold 2 : 准确率为: 0.978 耗时: 0.019946 s
Fold 3 : 准确率为: 0.978 耗时: 0.019946 s
Fold 4 : 准确率为: 0.978 耗时: 0.028446 s
Fold 5 : 准确率为: 0.977 耗时: 0.020291 s
平均准确率为: 0.978

```

平均准确率为97.8%，并且计算时间减少至 $10^{-2}$ 级，是比较可观的。

## 3. Part2

### 2.1 算法选择

分别选用BP单层神经网络、KNN、SVM、决策树、Bayes、多元回归、MLP、随机森林算法进行测试。

### 2.2 预处理

- 选取两队杀敌总数、死亡总数、助攻总数、经济和共四个特征，做归一化。
- 将两个csv文件合并，取后60000个用于训练与测试，前20000个来预测。
- 将时间标签化（除了多元回归）

```
1 def roundup(x):
2     return int(math.ceil(float(x) / 35)) * 35
```

## 2.3 算法与结果

### 2.3.1 BP单层神经网络

```
1  ##BP单层神经网络
2  import tensorflow.compat.v1 as tf
3  #定义参数
4  d=4 #输入节点个数
5  l=1 #输出节点个数
6  q=2*d+1 #隐层结点个数,采用经验公式2d+1
7  eta=0.5 #学习率
8  error=0.0016 #精度
9
10 #初始化权值和阈值
11 w1= tf.Variable(tf.random.normal([d, q], stddev=1, seed=1))
    #seed设定随机种子,保证每次初始化相同数据
12 b1=tf.Variable(tf.constant(0.0,shape=[q]))
13 w2= tf.Variable(tf.random.normal([q, l], stddev=1, seed=1))
14 b2=tf.Variable(tf.constant(0.0,shape=[l]))
15
16 #输入占位
17 tf.compat.v1.disable_eager_execution()
18 x = tf.placeholder(tf.float32, shape=(None, d))
19 y_ = tf.placeholder(tf.float32, shape=(None, l))
20
21 #前向传播
22 a=tf.nn.sigmoid(tf.matmul(x,w1)+b1) #sigmoid激活函数
23 y=tf.nn.sigmoid(tf.matmul(a,w2)+b2)
24 mse = tf.reduce_mean(tf.square(y_ - y)) #损失函数采用均方误差
25 train_step = tf.train.AdamOptimizer(eta).minimize(mse) #Adam
    算法
26
27 data = totalX
28 labels = totalY
29 labels = (labels - Min[4])/(Max[4]- Min[4])
30
31 #创建会话来执行图
```

```

32 with tf.Session() as sess:
33     init_op = tf.global_variables_initializer()#初始化节点
34     sess.run(init_op)
35
36     STEPS=0
37     while True:
38         sess.run(train_step, feed_dict={x: data, y_: labels})
39         STEPS+=1
40         train_mse= sess.run(mse, feed_dict={x: data, y_:
labels})
41         if STEPS % 10 == 0:#每训练100次，输出损失函数
42             print("第 %d 次训练后,训练集损失函数为: %g" % (STEPS,
train_mse))
43             if train_mse < error:
44                 break
45             print("总训练次数: ",STEPS)
46
47             #测试
48             Normal_y= sess.run(y, feed_dict={x: testX})#求得测试集下的y
计算值
49             DeNormal_y=Normal_y*(Max[4] - Min[4])+Min[4] #将y反归一化
50             Dy = DeNormal_y.tolist()
51             DY = testY.tolist()
52             n = len(Dy)
53             right = 0
54
55             for i in range(n):
56                 if Dy[i][0] >= (math.floor(DY[i][0]/100)*100-150) and
Dy[i][0] <= (math.ceil(DY[i][0]/100)*100+150):
57                     right += 1
58             print('正确率为: ',right/n)
59             err_BP = sess.run(mse, feed_dict={y: DeNormal_y, y_:
testY})#计算均方误差
60             print("测试集均方误差为: ",err_BP)
61
62             #预测
63             test = np.array(Test)
64             y= sess.run(y, feed_dict={x: test})
65             y = y*(Max[4] - Min[4])+Min[4] #将y反归一化

```

结果:

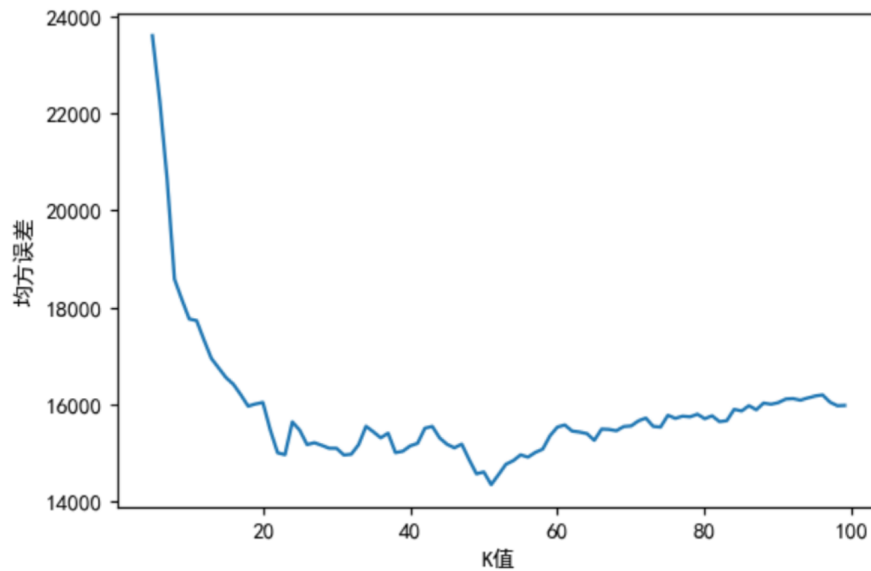


总训练次数: 478  
正确率为: 0.95175  
测试集均方误差为: 16664.203

### 2.3.2 KNN

```
1  ##KNN
2  from sklearn import neighbors
3  data = totalX
4  labels = totalY
5  E = [1000000,0]
6  err_KNN = [0] * 95
7
8  #遍历寻找最佳k值
9  for k in range(5,100):
10     knn = neighbors.KNeighborsClassifier(k)
11     knn.fit(data, labels)
12     y = knn.predict(testX)
13     #均方误差
14     err_KNN[k-5] = metrics.mean_squared_error (y, testY)
15     if err_KNN[k-5] < E[0]:
16         E[0] = err_KNN[k-5]
17         E[1] = k
18         Y = y
19
20 n = len(y)
21 right = 0
22 for i in range(n):
23     if Y[i] >= (math.floor(testY[i][0]/100)*100-150) and Y[i]
24     <= (math.ceil(testY[i][0]/100)*100+150):
25         #if abs(Y[i]-testY[i][0]) <= 100:
26         right += 1
27
28 plt.figure(dpi=110)
29 plt.xlabel("K值")
30 plt.ylabel("均方误差")
31 plt.plot(range(5,100),err_KNN)
32 plt.show()
33 print('最佳正确率为: ',right/n)
34 print("最佳K值为: ",E[1],"测试集均方误差为: ",E[0]);
```

结果:



最佳正确率为: 0.9745833333333334

最佳K值为: 51 测试集均方误差为: 14351.010333333334

### 2.3.3 SVM

```
1 #SVM
2 from sklearn import svm
3 data = totalX
4 labels = totalY
5 clf = svm.SVR()
6 clf.fit(totalX, totalY)
7 y = clf.predict(testX)
8 err_svm = metrics.mean_squared_error (y, testY)
9 print("测试集均方误差为: ",err_svm);
```

结果:

测试集均方误差为: 23766.188498855518

### 2.3.4 决策树

```

1  ##决策树
2  from sklearn import tree
3  data = totalX
4  labels = totalY
5  clf = tree.DecisionTreeRegressor()
6  clf = clf.fit(data, labels)
7  y = clf.predict(testX)
8  err_tree = metrics.mean_squared_error (y, testY)
9  print("测试集均方误差为: ",err_tree);

```

结果:

---

测试集均方误差为: 24150.702625

### 2.3.5 贝叶斯

```

1  #贝叶斯
2  from sklearn.naive_bayes import GaussianNB
3  data = totalX
4  labels = totalY
5  gnb = GaussianNB()
6  gnb = gnb.fit(data, labels)
7  y = gnb.predict(testX)
8  err_bayes = metrics.mean_squared_error (y, testY)
9  print("测试集均方误差为: ",err_bayes)

```

结果:

测试集均方误差为: 52772.974916666666

### 2.3.6 多元回归

```

1  #要重新预处理，时间不预处理（不使用roundup函数）
2  from sklearn.preprocessing import PolynomialFeatures
3  from sklearn.linear_model import LinearRegression
4
5  data = totalX
6  labels = totalY
7
8  reg = PolynomialFeatures(degree=4)
9  data_train = reg.fit_transform(data)

```

```
10 testX_train = reg.fit_transform(testX)
11 lg = LinearRegression()
12 lg.fit(data_train, labels)
13 y = lg.predict(testX_train)
14 err_ploy = metrics.mean_squared_error (y, testY)
15 print("测试集均方误差为: ", err_ploy)
```

结果:

测试集均方误差为: 12239.60481076169

### 2.3.7 MLP

```
1 from sklearn.neural_network import MLPClassifier
2
3 data = totalX
4 labels = totalY
5
6 clf = MLPClassifier(solver='lbfgs', alpha=1e-
7 5, hidden_layer_sizes=(3, 2), random_state=1)
8 clf.fit(data, labels)
9 y = clf.predict(testX)
10 err_mlp = metrics.mean_squared_error (y, testY)
11 print("测试集均方误差为: ", err_mlp)
```

结果:

测试集均方误差为: 25560.91075

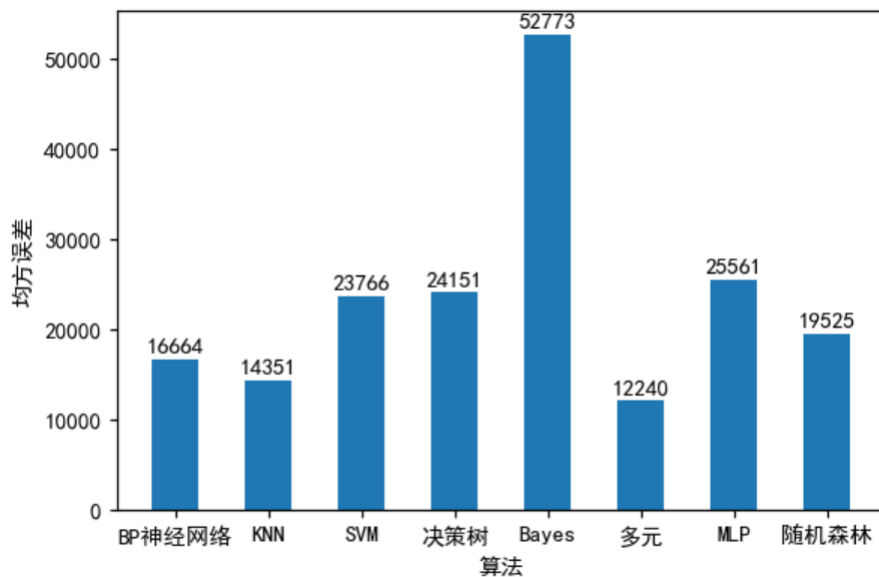
### 2.3.8 随机森林

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 data = totalX
4 labels = totalY
5
6 clf = RandomForestClassifier(n_jobs=3)
7 clf.fit(data, labels)
8 y = clf.predict(testX)
9 err_for = metrics.mean_squared_error (y, testY)
10 print("测试集均方误差为: ", err_for)
```

结果:

测试集均方误差为: 19525.337083333332

## 2.4 结果对比



在系统训练集、测试集下，均方误差最小的为**KNN算法与多元回归算法**。

## 2.5 Autogluon库

```
1 from autogluon.tabular import TabularDataset, TabularPredictor
2 predictor =
  TabularPredictor(label='gameDuration').fit(df[0:48000])
3 a=df[48000:60000].drop(['gameDuration'],axis=1)
4 y = predictor.predict(a)
5 err_auto = metrics.mean_squared_error (y, testY)
6 print("测试集均方误差为: ",err_auto);
```

结果:

测试集均方误差为: 11266.0005795075

准确度最佳。

## 参考文献

- [1] 宋毅飞,周剑秋.KNN算法与其改进算法的性能比较[J].机电产品开发与创新,2017,30(02):60-63.