# 模式识别与机器学习课程作业

| | |
|---|---|
| 学院 | 中国科学院大学 人工智能学院 |
| 姓名 | 芮志清 |
| 学号 | 2018Z8020661080 |
| 邮箱 | ruizhiqing18🖳mails.ucas.ac.cn |
| 时间 | 2019年5月31日 |

## 作业要求

老师给一些数据，从中任选一个，自选主题，做一个实验

截止日期：**6月1日**

## 选题

### 数据集

本次选择的数据为"Bank Marketing Data"[1] （银行市场数据）。该数据与葡萄牙银行机构的电话推广活动有关。 数据包含该数据集的分类目标为预测客户是否订阅定期存款服务。

本次实验采用数据集的标准子集，包含以下特征列：

1. 年龄
2. 工作类型
3. 婚姻状况
4. 教育
5. 默认信用
6. 是否有住房贷款
7. 是否有个人贷款
8. 联系人沟通方式
9. 上次联系的在每月的第几天
10. 上次联系的月份
11. 上次联系为周几
12. 沟通时长
13. 与客户的联系次数
14. 上一次推广后经过的天数
15. 此活动前和此客户之前的联系次数
16. 上一次营销活动的结果
17. 有客户是否订阅定期存款?

## 实验目标

搭建神经网络模型，实现对用户订阅存款的预测。

## 实验分析

### 数据预处理

本次数据集为结构化数据，包括17列，其中16列为特征列，1列为标签列。特征列中包含连续值、类别和布尔值，但是类别和布尔值均为字符串表示，需要将其进行数值化处理。

标签为布尔值，因此本实验为一个**二分类**的预测问题。

在类别特征中，类别之间并无数值大小的关系，因此对其进行onehot编码。

在onehot编码后，出现了较多特征，因此可以计算特征和目标的权重，取一部分特征，进行特征降维，准备使用xgboost实现此功能。

在特征降维之前，为了使得xgboost使用更加准确，因此在onehot编码后先进行归一化处理。

### 搭建模型选择

本次实验准备采用tensorflow及keras搭建多层全连接神经网络模型。

在层间采用**relu**激活函数，在输出层采用**sigmod**激活函数，优化器采用**rmsprop**，损失函数采用**binary_crossentropy**。

# 实验进行情况

## 实验环境

本次实验采用python3作为编程语言，实验代码运行在Google Colab上([ML-Assignment](#))，代码工程托管在github上([ZQRui / ML-Assignment](#))

## 代码及运行情况

### 使用的库

```
1  import pandas as pd
2  import numpy as np
3  import tensorflow as tf
4  from tensorflow import keras
5  from tensorflow.keras import layers
6  import matplotlib.pyplot as plt
7  import xgboost as xgb
8  import json
9  import operator
10 from collections import OrderedDict
11 import os,sys
12 import logging
13
```

### 设置日志输出

```
1  #DEBUG to const.LOGFILE
2  logging.basicConfig(level=logging.DEBUG,
3                      format='%(asctime)s %(filename)s[line:%(lineno)d] %(levelname)s %(message)s', # 输出
   格式
4                      datefmt='%a, %d %b %Y %H:%M:%S',
5                      filename="debug.log.txt",
6                      filemode='a')
7
8  # to console
9  console = logging.StreamHandler()
10 console.setLevel(logging.INFO)
11 formatter = logging.Formatter('%(asctime)s: %(levelname)-5s %(message)s')
12 console.setFormatter(formatter)
13 logging.getLogger('').addHandler(console)
14
15 #INFO to const.INFOLOGFILE
16 infolog = logging.FileHandler("info.log.txt")
17 infolog.setLevel(logging.INFO)
18 errformatter = logging.Formatter('%(asctime)s %(filename)s[line:%(lineno)d] %(levelname)s %(message)s')
19 infolog.setFormatter(errformatter)
20 logging.getLogger('').addHandler(infolog)
21
22 #ERROR to const.ERRLOGFILE
23 errlog = logging.FileHandler("error.log.txt")
24 errlog.setLevel(logging.WARNING)
25 errformatter = logging.Formatter('%(asctime)s %(filename)s[line:%(lineno)d] %(levelname)s %(message)s')
26 errlog.setFormatter(errformatter)
27 logging.getLogger('').addHandler(errlog)
28
```

### 读数据

```
1  dataset_file="./bank.csv"
```

```
1  df=pd.read_csv(dataset_file,sep=";")
2  df.head()
```

| | age | job | marital | education | default | balance | housing | loan | contact | day | month | duration | campa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 30 | unemployed | married | primary | no | 1787 | no | no | cellular | 19 | oct | 79 | 1 |
| **1** | 33 | services | married | secondary | no | 4789 | yes | yes | cellular | 11 | may | 220 | 1 |
| **2** | 35 | management | single | tertiary | no | 1350 | yes | no | cellular | 16 | apr | 185 | 1 |
| **3** | 30 | management | married | tertiary | no | 1476 | yes | yes | unknown | 3 | jun | 199 | 4 |
| **4** | 59 | blue-collar | married | secondary | no | 0 | yes | no | unknown | 5 | may | 226 | 1 |

**数据集的描述**

```
1 df.describe()
```

| | age | balance | day | duration | campaign | pdays | previous |
|---|---|---|---|---|---|---|---|
| **count** | 4521.000000 | 4521.000000 | 4521.000000 | 4521.000000 | 4521.000000 | 4521.000000 | 4521.000000 |
| **mean** | 41.170095 | 1422.657819 | 15.915284 | 263.961292 | 2.793630 | 39.766645 | 0.542579 |
| **std** | 10.576211 | 3009.638142 | 8.247667 | 259.856633 | 3.109807 | 100.121124 | 1.693562 |
| **min** | 19.000000 | -3313.000000 | 1.000000 | 4.000000 | 1.000000 | -1.000000 | 0.000000 |
| **25%** | 33.000000 | 69.000000 | 9.000000 | 104.000000 | 1.000000 | -1.000000 | 0.000000 |
| **50%** | 39.000000 | 444.000000 | 16.000000 | 185.000000 | 2.000000 | -1.000000 | 0.000000 |
| **75%** | 49.000000 | 1480.000000 | 21.000000 | 329.000000 | 3.000000 | -1.000000 | 0.000000 |
| **max** | 87.000000 | 71188.000000 | 31.000000 | 3025.000000 | 50.000000 | 871.000000 | 25.000000 |

# 建立onehot映射表

```
1 onehot_cols=["job","marital","education","default","housing","loan","contact","month","poutcome"]
```

```
1 for col in onehot_cols:
2   df[col]=df[col].astype("category")
```

```
1 onehot_newkey=[]
2 for col in onehot_cols:
3   for v in df[col].unique():
4     c=(f"{col}-{v}",col,v)
5     onehot_newkey.append(c)
```

```
1 onehot_newkey.append(("pdays--1","pdays",-1))
2 onehot_newkey.append(("pdays-yes","pdays--1",0))
```

## Onehot编码后数据列的映射表

```
1 onehot_newkey
```

```
1  # 分别对应编码后的列名，编码前的列名，编码值
2  [('job-unemployed', 'job', 'unemployed'),
3   ('job-services', 'job', 'services'),
4   ('job-management', 'job', 'management'),
5   ('job-blue-collar', 'job', 'blue-collar'),
6   ('job-self-employed', 'job', 'self-employed'),
7   ('job-technician', 'job', 'technician'),
8   ('job-entrepreneur', 'job', 'entrepreneur'),
9   ('job-admin.', 'job', 'admin.'),
10  ('job-student', 'job', 'student'),
11  ('job-housemaid', 'job', 'housemaid'),
12  ('job-retired', 'job', 'retired'),
13  ('job-unknown', 'job', 'unknown'),
14  ('marital-married', 'marital', 'married'),
15  ('marital-single', 'marital', 'single'),
16  ('marital-divorced', 'marital', 'divorced'),
17  ('education-primary', 'education', 'primary'),
18  ('education-secondary', 'education', 'secondary'),
19  ('education-tertiary', 'education', 'tertiary'),
```

```
20    ('education-unknown', 'education', 'unknown'),
21    ('default-no', 'default', 'no'),
22    ('default-yes', 'default', 'yes'),
23    ('housing-no', 'housing', 'no'),
24    ('housing-yes', 'housing', 'yes'),
25    ('loan-no', 'loan', 'no'),
26    ('loan-yes', 'loan', 'yes'),
27    ('contact-cellular', 'contact', 'cellular'),
28    ('contact-unknown', 'contact', 'unknown'),
29    ('contact-telephone', 'contact', 'telephone'),
30    ('month-oct', 'month', 'oct'),
31    ('month-may', 'month', 'may'),
32    ('month-apr', 'month', 'apr'),
33    ('month-jun', 'month', 'jun'),
34    ('month-feb', 'month', 'feb'),
35    ('month-aug', 'month', 'aug'),
36    ('month-jan', 'month', 'jan'),
37    ('month-jul', 'month', 'jul'),
38    ('month-nov', 'month', 'nov'),
39    ('month-sep', 'month', 'sep'),
40    ('month-mar', 'month', 'mar'),
41    ('month-dec', 'month', 'dec'),
42    ('poutcome-unknown', 'poutcome', 'unknown'),
43    ('poutcome-failure', 'poutcome', 'failure'),
44    ('poutcome-other', 'poutcome', 'other'),
45    ('poutcome-success', 'poutcome', 'success'),
46    ('pdays--1', 'pdays', -1),
47    ('pdays-yes', 'pdays--1', 0)]
```

```
1  conf={"onehot_cols":onehot_cols,"onehot_newkey":onehot_newkey}
```

```
1  with open("conf.json","w") as fp:
2      json.dump(conf,fp,indent=2)
```

## 对数据进行Onehot编码及归一化处理

```
1  df["y"]=(df["y"]=="yes").astype(float)
```

```
1  def onehot(df,newkeys,oldkeys):
2      df=df.copy()
3      for key,oldkey,value in newkeys:
4          logging.debug(f"{key},{oldkey},{value}")
5          df[key]=(df[oldkey]==value).astype(float)
6      for oldkey in oldkeys:
7          df.pop(oldkey)
8      return df
```

```
1  df2=onehot(df,onehot_newkey,onehot_cols)
```

```
1  def split_train_test(source, frac, ):
2      train_dataset = source.sample(frac=0.8, random_state=0)
3      test_dataset = source.drop(train_dataset.index)
4
5      return train_dataset, test_dataset
```

```
1  train_dataset,test_dataset=split_train_test(df2,0.8)
```

```
1  # y=train_dataset.pop("y")
2  train_labels = train_dataset.pop('y')
3  test_labels = test_dataset.pop('y')
4
5
```

```
1  stats=train_dataset.describe().transpose()
2  print(stats)
```

```
1                  count         mean          std  ...    50%     75%      max
2  age            3617.0    41.115842    10.573495  ...   39.0    49.0     87.0
3  balance        3617.0  1405.922311  2972.465627  ...  444.0  1465.0  71188.0
4  day            3617.0    15.848217     8.220174  ...   16.0    21.0     31.0
5  duration       3617.0   268.094001   265.199283  ...  187.0   333.0   3025.0
```

```
6   campaign            3617.0    2.809511    3.137596  ...    2.0    3.0   50.0
7   pdays               3617.0   39.948023  100.672342  ...   -1.0   -1.0  871.0
8   previous            3617.0    0.553221    1.729015  ...    0.0    0.0   25.0
9   job-unemployed      3617.0    0.029030    0.167913  ...    0.0    0.0    1.0
10  job-services        3617.0    0.093171    0.290712  ...    0.0    0.0    1.0
11  job-management      3617.0    0.215648    0.411328  ...    0.0    0.0    1.0
12  job-blue-collar     3617.0    0.210395    0.407646  ...    0.0    0.0    1.0
13  job-self-employed   3617.0    0.040918    0.198127  ...    0.0    0.0    1.0
14  job-technician      3617.0    0.167819    0.373757  ...    0.0    0.0    1.0
15  job-entrepreneur    3617.0    0.037600    0.190254  ...    0.0    0.0    1.0
16  job-admin.          3617.0    0.107271    0.309501  ...    0.0    0.0    1.0
17  job-student         3617.0    0.017971    0.132863  ...    0.0    0.0    1.0
18  job-housemaid       3617.0    0.024606    0.154943  ...    0.0    0.0    1.0
19  job-retired         3617.0    0.048659    0.215184  ...    0.0    0.0    1.0
20  job-unknown         3617.0    0.006912    0.082861  ...    0.0    0.0    1.0
21  marital-married     3617.0    0.613215    0.487081  ...    1.0    1.0    1.0
22  marital-single      3617.0    0.264860    0.441320  ...    0.0    1.0    1.0
23  marital-divorced    3617.0    0.121924    0.327244  ...    0.0    0.0    1.0
24  education-primary   3617.0    0.150401    0.357513  ...    0.0    0.0    1.0
25  education-secondary 3617.0    0.508432    0.499998  ...    1.0    1.0    1.0
26  education-tertiary  3617.0    0.299419    0.458067  ...    0.0    1.0    1.0
27  education-unknown   3617.0    0.041747    0.200039  ...    0.0    0.0    1.0
28  default-no          3617.0    0.984241    0.124559  ...    1.0    1.0    1.0
29  default-yes         3617.0    0.015759    0.124559  ...    0.0    0.0    1.0
30  housing-no          3617.0    0.424938    0.494402  ...    0.0    1.0    1.0
31  housing-yes         3617.0    0.575062    0.494402  ...    1.0    1.0    1.0
32  loan-no             3617.0    0.851811    0.355336  ...    1.0    1.0    1.0
33  loan-yes            3617.0    0.148189    0.355336  ...    0.0    0.0    1.0
34  contact-cellular    3617.0    0.644457    0.478744  ...    1.0    1.0    1.0
35  contact-unknown     3617.0    0.291125    0.454344  ...    0.0    1.0    1.0
36  contact-telephone   3617.0    0.064418    0.245530  ...    0.0    0.0    1.0
37  month-oct           3617.0    0.016865    0.128783  ...    0.0    0.0    1.0
38  month-may           3617.0    0.311861    0.463317  ...    0.0    1.0    1.0
39  month-apr           3617.0    0.060271    0.238021  ...    0.0    0.0    1.0
40  month-jun           3617.0    0.118883    0.323696  ...    0.0    0.0    1.0
41  month-feb           3617.0    0.050041    0.218061  ...    0.0    0.0    1.0
42  month-aug           3617.0    0.140448    0.347499  ...    0.0    0.0    1.0
43  month-jan           3617.0    0.031794    0.175476  ...    0.0    0.0    1.0
44  month-jul           3617.0    0.156483    0.363363  ...    0.0    0.0    1.0
45  month-nov           3617.0    0.086812    0.281599  ...    0.0    0.0    1.0
46  month-sep           3617.0    0.010782    0.103291  ...    0.0    0.0    1.0
47  month-mar           3617.0    0.011059    0.104593  ...    0.0    0.0    1.0
48  month-dec           3617.0    0.004700    0.068405  ...    0.0    0.0    1.0
49  poutcome-unknown    3617.0    0.818911    0.385145  ...    1.0    1.0    1.0
50  poutcome-failure    3617.0    0.110036    0.312978  ...    0.0    0.0    1.0
51  poutcome-other      3617.0    0.043406    0.203798  ...    0.0    0.0    1.0
52  poutcome-success    3617.0    0.027647    0.163983  ...    0.0    0.0    1.0
53  pdays--1            3617.0    0.818911    0.385145  ...    1.0    1.0    1.0
54  pdays-yes           3617.0    0.181089    0.385145  ...    0.0    0.0    1.0
55
56  [53 rows x 8 columns]
```

```
1  def norm(df,stats):
2    return (df - stats['mean'])/stats['std']
```

```
1  normed_train_data=norm(train_dataset,stats)
2  normed_test_data = norm(test_dataset, stats)
```

## 特征权重计算

```
1   def feature_importance(df,xl,yl):
2       df=df.copy()
3       features=xl
4       with open('xgb.fmap',"w") as fpmap:
5           i=0
6           for fe in features:
7               fpmap.write(f"{i}\t{fe}\tq\n")
8               i=i+1
9       params = {
10              'min_child_weight': 0,
11              'eta': 0.02,
12              'colsample_bytree': 0.7,
13              'max_depth': 12,
14              'subsample': 0.7,
15              'alpha': 1,
16              'gamma': 1,
```

```
17              'silent': 1,
18              'verbose_eval': True,
19              'seed': 12
20          }
21
22      rounds=100
23      y=df[yl]
24
25      X=df[Xl]
26
27      xgtrain=xgb.DMatrix(X,label=y)
28      bst=xgb.train(params,xgtrain,num_boost_round=rounds)
29      importance=bst.get_fscore(fmap="xgb.fmap")
30
31      importance=sorted(importance.items(),key=operator.itemgetter(1),reverse=True)
32
33      return importance
```

```
1   all_cols=list(df.columns)
2   all_cols.remove("y")
3   feature_sel_data=normed_train_dataset.copy()
4   feature_sel_data["y"]=train_labels
5   importance=feature_importance(feature_sel_data,normed_train_dataset.columns,"y")
6   importance
```

**特征权重值:**

```
1   [('duration', 275),
2    ('age', 131),
3    ('day', 113),
4    ('pdays', 87),
5    ('balance', 80),
6    ('poutcome-success', 71),
7    ('month-oct', 51),
8    ('month-mar', 49),
9    ('previous', 39),
10   ('contact-unknown', 30),
11   ('month-apr', 28),
12   ('housing-no', 25),
13   ('marital-married', 23),
14   ('month-feb', 22),
15   ('education-tertiary', 21),
16   ('campaign', 20),
17   ('month-jun', 17),
18   ('poutcome-other', 13),
19   ('contact-cellular', 13),
20   ('month-may', 13),
21   ('month-nov', 10),
22   ('month-sep', 9),
23   ('contact-telephone', 7),
24   ('month-dec', 7),
25   ('loan-no', 7),
26   ('job-management', 6),
27   ('job-blue-collar', 6),
28   ('job-student', 6),
29   ('month-jul', 5),
30   ('job-retired', 4),
31   ('job-technician', 4),
32   ('housing-yes', 4),
33   ('month-aug', 4),
34   ('poutcome-failure', 4),
35   ('default-yes', 3),
36   ('default-no', 3),
37   ('month-jan', 3),
38   ('job-entrepreneur', 2),
39   ('job-housemaid', 2),
40   ('education-primary', 2),
41   ('job-unemployed', 2),
42   ('marital-divorced', 2),
43   ('education-unknown', 2),
44   ('job-unknown', 2),
45   ('pdays--1', 1),
46   ('loan-yes', 1)]
```

## 神经网络模型

**超参数**

```
1  args={"feature_count": 6, "train_data_frac": 0.85, "layers_count": [6, 5], "epochs": 5000}
2
```

**取部分特征**

```
1  import_features = [feature for feature, v in importance[:feature_count]]
2  normed_train_data=normed_train_data[import_features]
3  normed_test_data = normed_test_data[import_features]
```

**神经网络模型**

```
1  def build_model():
2      layers_list = [layers.Dense(layers_count[0], activation=tf.nn.relu, input_shape=[feature_count]), ] + \
3                     [layers.Dense(count, activation=tf.nn.relu, ) for count in layers_count[1:]] + \
4                     [layers.Dense(1, activation=tf.nn.sigmoid)]
5      model = keras.Sequential(layers_list)
6      model.compile(optimizer='rmsprop',loss='binary_crossentropy',metrics=['accuracy'])
7
8      return model
```

**模型训练**

```
1   logging.info(f"evaluate model with arg {args}")
2   feature_count = args.get("feature_count", 5)
3
4   layers_count = args.get("layers_count", [5])
5   logging.debug(train_dataset.describe())
6
7   model=build_model()
8   EPOCHS = args.get("epochs", 1000)
9   early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', patience=10)
10  logging.debug(normed_train_data.describe())
11  history = model.fit(
12      normed_train_data, train_labels,
13      epochs=EPOCHS, validation_split=0.2, verbose=0, callbacks=[early_stop]
14  )
15  test_loss, test_acc = model.evaluate(normed_test_data, test_labels)
16
17  logging.info(f"acc: {test_acc} loss:{test_loss} args: {args}")
```

**训练结果**

```
1   2019-05-30 22:31:58,507: INFO  evaluate model with arg {'feature_count': 6, 'train_data_frac': 0.85,
    'layers_count': [6, 5], 'epochs': 5000, }
2
3
4   WARNING:tensorflow:From /usr/local/lib/python3.6/dist-
    packages/tensorflow/python/ops/math_ops.py:3066: to_int32 (from tensorflow.python.ops.math_ops) is
    deprecated and will be removed in a future version.
5   Instructions for updating:
6   Use tf.cast instead.
7
8
9   2019-05-30 22:31:58,812: WARNING From /usr/local/lib/python3.6/dist-
    packages/tensorflow/python/ops/math_ops.py:3066: to_int32 (from tensorflow.python.ops.math_ops) is
    deprecated and will be removed in a future version.
10  Instructions for updating:
11  Use tf.cast instead.
12
13
14  904/904 [==============================] - 0s 30us/sample - loss: 0.2311 - acc: 0.9126
15
```

**模型的测试准确率**

模型的准确率达到了91%

```
1  2019-05-30 22:32:03,508: INFO  acc: 0.9126105904579163 loss:0.23113595727270683 args: {'feature_count': 6,
   'train_data_frac': 0.85, 'layers_count': [6, 5], 'epochs': 5000, }
2
```

**超参数优化**

本实验还采用了超参数的优化，采用了随机方法选择参数。

```
1   def model_evaluate_times(d, times):
2       s = 0
3       for i in range(times):
4           logging.debug(f"evaluate times : {times}")
5           s += model_evaluate(d)
6       s /= times
7       logging.info(f"model acc average = {s}")
8       return s
9
10
11  def generate_args():
12      l = []
13      for epochs in [100, 500, 1000, 5000, 10000]:
14          for frac in [0.6, 0.7, 0.75, 0.8, 0.85, 0.9]:
15              for layers_count in [[i] for i in range(1, 10)] + \
16                                  [[i, j] for i in range(1, 10) for j in range(1, 10)]:
17              #   [[i, j, k] for i in range(1, 10) for j in range(1, 10) for k in range(1, 10)]:
18                  for fc in range(2, 10):
19                      d = {"feature_count": fc, "train_data_frac": frac, "layers_count": layers_count,
    "epochs": epochs}
20                      l.append(d)
21                      logging.debug(f"generate args {d}")
22                      yield d
```

# 实验总结

## 收获

- 通过本实验，学习了使用TensorFlow搭建多层神经网络模型解决二分类问题
- 采用了OneHot编码方案
- 使用了随机超参数优化方案
- 达到了91%的准确率

## 问题及解决方案

- 由于超参数优化比较消耗时间，因此不容易寻找到准确率更高的方案。
    - 采用随机方法寻找在超参数空间寻找。
    - 对于数值型超参数，从粗度较大的参数表开始，根据已得到的性能较高的超参数为基础，再进行细粒度的超参数搜索。

---

1. In P. Novais et al. (Eds.), Proceedings of the European Simulation and Modelling Conference - ESM'2011, pp. 117-121, Guimarães, Portugal, October, 2011. EUROSIS.↩