

ANN Classifier Assignment Assignment 3

Xiao(Megan)Ling

1. Object

We are going to use UCI red wine dataset predict white wine quality using an Artificial Neural Network. With given data points about Portuguese “Vinho Verde” wine, such as fixed acidity, residual sugar, etc., we will be able to determine wine qualities based on a model designed.

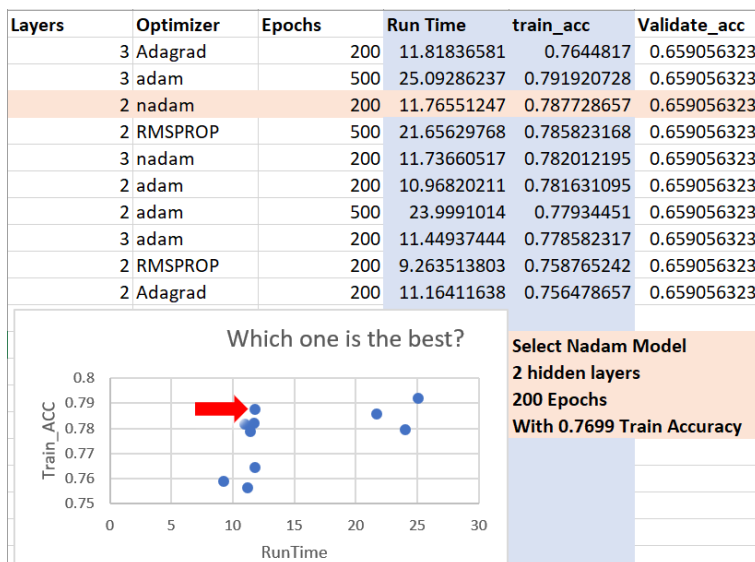
2. The Final Model & Training Algorithm

Final Model Application

```
model = Sequential()
model.add(Dense(12, activation='relu', input_shape=(x_train.shape[1],)))
model.add(Dense(12, activation='relu'))
model.add(Dense(1, activation = 'sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='Nadam', metrics=['accuracy'])
new = model.fit(x_test, y_test, epochs=80, batch_size=16, verbose=0)
train_acc_new = new.history["acc"]
print(max(train_acc_new))
```

The Final Model to predict the wine quality is a binary classification model, which has 2 hidden layers using Nadam as the model optimizer, binary_crossentropy as loss function and Accuracy as metrics. The training algorithm here is Artificial Neural Network Algorithms which consists of different layers with weighted connections for analyzing and learning data.

3. The experimental plan for arriving at the final model and How long it took to run all the models in your experimental plan



To arrive at the final model, I changed parameters such as the number of Epochs, the size of the hidden layers, and the model compiler optimizers.

First, I designed model with RMSprop as compile optimizers with Epochs = 500 and 2 hidden layers. The RMSprop is suitable selection since it divides the learning rate by exponential decaying average of squared gradients. I applied the same parameters on different model optimizers, Nadam and Adam. I chose Adam because it requires low memory and Nadam is Adam RMSprop with Nesterov

momentum. Also, I include Adagrad. It adapts the learning rate to the parameters, performing smaller updates for parameters associated with frequently occurring features.

The considered time factor, generally, models with 500 epochs run slower than models with 200 epochs but with higher training accuracy. However, some models with 200 epochs and additional hidden layers can reach similar level training accuracy and save time. As a result, I selected the Nadam Model with 2 hidden layers and 200 epochs as my final model.

4. An explanation of the input variables and any preprocessing steps you took

The screenshot shows an Excel spreadsheet titled 'winequality-white.csv - Excel'. The data is organized into columns: chlorides, free sulfur, total sulfur, density, pH, sulphates, alcohol, and quality. A sorting filter menu is open over the 'quality' column, showing options to sort by color, filter by color, and number filters. The 'Number Filters' section is expanded, showing a search bar and a list of values: (Select All), 0, and 1. The 'OK' button is highlighted.

	E	F	G	H	I	J	K	L
	chlorides	free sulfur	total sulfur	density	pH	sulphates	alcohol	quality
1	0.036	20	114	0.99248	3.75			
2	0.03	29	118	0.989	3.57			
8	0.029	93	161	0.98999	3.65			
1	0.041	64	157	0.99688	3.42			
8	0.024	31	111	0.98816	3.48			
3	0.03	31	127	0.98904	3.46			
1	0.038	52	97	0.99022	3.41			
5	0.033	89	159	0.99332	3.34			
4	0.053	11	178	0.99426	3.79			
1	0.02	5	9	0.98722	3.3			
9	0.036	33	106	0.98746	3.21			
1	0.042	35	90	0.9908	3.76			
3	0.036	69	168	0.99212	3.47			
4	0.036	23	134	0.98981	3.53			
5	0.028	33	163	0.9939	3.36			
5	0.028	34	163	0.9937	3.35			
2	0.037	17	112	0.99324	3.66			
9	0.03	22	111	0.9902	3.38			
6	0.034	23	111	0.99274	3.46			
1	0.013	4	10	0.99246	3.32			
2	0.042	40	98	0.9898	3.42			
2	0.074	47	130	0.99132	3.31			
1	0.044	28	100	0.99034	3.38			

To figure out the propitiate separate point for binary classification, I calculated the mean of the qualities, and it returns around 5.6. Consequently, I set 0 – 5 wine quality as low quality and 6 – 10 as high wine quality.

The challenge of binary classification is to put data into right categories. I tried to use `pd.replace` to change target variables Y into two categories 0 and 1, which means low and high qualities. But it costs me a long time to figure out the Python function. Then I realized since the white wine dataset is a relatively small data set, it would nice and quick to use Excel manually categorize qualities into 0 and 1.

Pretreat Data Section

```
x_data = data[x_vars]
y_data = data[y_vars]

#print (len(x_data))
#print (len(y_data))

x_data_MinMax = preprocessing.MinMaxScaler()

x_data.values
x_data = np.array(x_data).reshape((len(x_data), 11))

y_data.values
y_data = np.array(y_data).reshape((len(y_data), 1))

x_data = x_data_MinMax.fit_transform(x_data)

x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.33, random_
x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_

x_train.mean(axis=0)
y_train.mean(axis=0)
x_test.mean(axis=0)
y_test.mean(axis=0)
```

I also scaled variables x_data, which contains 11 variables such as chlorides, free sulfur dioxide, total sulfur dioxide etc. to the same measures. For target variable y_data set, I already changed them into binary values. So, there is no need to scale target variable dataset. Then splitting the dataset into train, test and validate datasets through train_test_split function.

6. An explanation of your model, metrix and justification for your choice

```
model = Sequential()
model.add(Dense(12, activation='relu', input_shape=(x_train.shape[1],)))
model.add(Dense(12, activation='relu'))
model.add(Dense(1, activation = 'sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='Adagrad', metrics=['accuracy'])
```

To prevent overfitting, I selected a very small network with three hidden layers, each with 12 units. The network ends with a single unit and no activation. Since we are building up a binary classifier model, binary_crossentropy loss function is widely applied when compiling a binary classifier network. Also, I add activation='sigmoid' in to the last model layer since 'sigmoid' indicates binary results. To monitor the training result, accuracy metric is introduced to show the difference between the predictions and the targets.

```
history = model.fit(x_train, y_train,
                    validation_data=(x_val, y_val),
                    epochs=200, batch_size=50, verbose=0)

train_loss = history.history["loss"]
validate_loss = history.history["val_loss"]
train_acc = history.history["acc"]
validate_acc = history.history["val_acc"]
```

For binary classification, it is prevalent to monitor binary_crossentropy loss and accuracy. As expected, the binary classification model accuracy is not as high as categorical_crossentropy. Binary_crossentropy applied when the data has two target classes, in this case are 0 and 1. The accuracy monitors if the model predict accurate results as the y from the test.set.

With the best model selected, nearly 80% of model prediction are correct. It is not a bad result compared with the complexity of building a 10 class multi-classification and bearing the risks of

overfitting. Plus, when we digging into the Excel data, we found out the lowest quality level of the Portuguese white wine is 3, which means it is unnecessary to build a 10 class model which has 3 empty class to fit. The multi-classification model prediction result would probably be a overfitting result.

fixed ac	volatile	citric ac	residua	chloride	free sul	total su	density	pH	sulphat	alcohol	quality
7	0.27	0.36	20.7	0.045	45	170	1.001	1.001			
6.3	0.3	0.34	1.6	0.049	14	132	0.994	0.994			
8.1	0.28	0.4	6.9	0.05	30	97	0.995	0.995			
7.2	0.23	0.32	8.5	0.058	47	186	0.995	0.995			
7.2	0.23	0.32	8.5	0.058	47	186	0.995	0.995			
8.1	0.28	0.4	6.9	0.05	30	97	0.995	0.995			
6.2	0.32	0.16	7	0.045	30	136	0.994	0.994			
7	0.27	0.36	20.7	0.045	45	170	1.001	1.001			
6.3	0.3	0.34	1.6	0.049	14	132	0.994	0.994			
8.1	0.22	0.43	1.5	0.044	28	129	0.993	0.993			
8.1	0.27	0.41	1.45	0.033	11	63	0.990	0.990			
8.6	0.23	0.4	4.2	0.035	17	109	0.994	0.994			
7.9	0.18	0.37	1.2	0.04	16	75	0.992	0.992			
6.6	0.16	0.4	1.5	0.044	48	143	0.991	0.991			
8.3	0.42	0.62	19.25	0.04	41	172	1.000	1.000			
6.6	0.17	0.38	1.5	0.032	28	112	0.991	0.991			
6.3	0.48	0.04	1.1	0.046	30	99	0.992	0.992			
6.2	0.66	0.48	1.2	0.029	29	75	0.989	0.989			
7.4	0.34	0.42	1.1	0.033	17	171	0.991	0.991			
6.5	0.31	0.14	7.5	0.044	34	133	0.995	0.995			

7. Business Insights

Even though binary classification model sacrificed some accuracy in prediction, from a practical perspective, I believe a binary classification model is a good fit for people who wants to invest in the Portuguese “Vinho Verde” wine business. As mentioned before, the white wine is seldomly categorized as a quality 1 level wine, which means the white wine producer would not put the actual low quality wine into the market. Plus, the wine quatli system within EU indicates that tasting wine quality should within the “high” level boundy. So, the white wine quality around 5-6 will probably sale in open market instead of in auction or private sales. Plus, based on demographic data and common sense, people who consume white wine has higher education, better knowledge of tasting, and higher income. Consumers and investors would care more about those wines who graded as good quality. A quick model would help them know their wine qualities within seconds.

CODE

Import Libraries Section

```
from keras.models import Sequential
from keras.layers import Dense
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
import numpy as np
import pandas as pd
import time
```

```
import sys
if "C:\\My_Python_Lib" not in sys.path:
    sys.path.append("C:\\My_Python_Lib")
```

Load Data Section

```
start = time.time()

data = pd.read_csv("winequality-white.csv",
                   header=0,
                   names = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
                           'chlorides', 'free sulfur dioxide',
                           'total sulfur dioxide', 'density', 'pH',
                           'sulphates', 'alcohol', 'quality'])

x_vars = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
          'chlorides', 'free sulfur dioxide',
          'total sulfur dioxide', 'density', 'pH',
          'sulphates', 'alcohol',]

y_vars = ['quality']
```

```
.....
```

Pretreat Data Section

```
.....
```

```
x_data = data[x_vars]
y_data = data[y_vars]
```

```
#print (len(x_data))
#print (len(y_data))
```

```
x_data_MinMax = preprocessing.MinMaxScaler()
```

```
x_data.values
x_data = np.array(x_data).reshape((len(x_data), 11))
```

```
y_data.values
y_data = np.array(y_data).reshape((len(y_data), 1))
```

```
x_data = x_data_MinMax.fit_transform(x_data)
```

```
x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.33, random_state=7)
x_train, x_val, y_train, y_val = train_test_split(x_train, y_train, test_size=0.2, random_state=7)
```

```
x_train.mean(axis=0)
y_train.mean(axis=0)
x_test.mean(axis=0)
y_test.mean(axis=0)
```

```
.....
```

Define Model Section - RMSPROP

```
.....
```

```
model = Sequential()
model.add(Dense(12, activation='relu', input_shape=(x_train.shape[1],)))
model.add(Dense(12, activation='relu'))
```

```
model.add(Dense(1,activation = 'sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='Adagrad', metrics=['accuracy'])
```

.....

ACC & LOSS validation

.....

```
history = model.fit(x_train, y_train,
                    validation_data=(x_val, y_val),
                    epochs=200, batch_size=50, verbose=0)
```

```
train_loss = history.history["loss"]
validate_loss = history.history["val_loss"]
train_acc = history.history["acc"]
validate_acc = history.history["val_acc"]
end = time.time()
print(end - start)
print(max(train_acc))
print(min(validate_acc))
```

.....

Final Model Application

.....

```
model = Sequential()
model.add(Dense(12, activation='relu', input_shape=(x_train.shape[1],)))
model.add(Dense(12, activation='relu'))
model.add(Dense(1,activation = 'sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='Nadam', metrics=['accuracy'])
new = model.fit(x_test, y_test, epochs=80, batch_size=16, verbose=0)
train_acc_new = new.history["acc"]
print(max(train_acc_new))
```