

PROGRAMMING ASSIGNMENT #2

T81-559: Applications of Deep Neural Networks, Washington University September 16, 2017

Listing 1 shows a sample submission skeleton that you can use as a starting point for this assignment.

Listing 1: Sample Submission Skeleton

```
1
2
3 path = "./data/"
4
5 # Helpful functions
6
7 # Encode text values to dummy variables(i.e. [1,0,0],[0,1,0],[0,0,1] for ↵
   red,green,blue)
8 def encode_text_dummy(df, name):
9     dummies = pd.get_dummies(df[name])
10    for x in dummies.columns:
11        dummy_name = "{}-{}".format(name, x)
12        df[dummy_name] = dummies[x]
13        df.drop(name, axis=1, inplace=True)
14
15
16 # Encode text values to a single dummy variable. The new columns (which ↵
   do not replace the old) will have a 1
17 # at every location where the original column (name) matches each of the ↵
   target_values. One column is added for
18 # each target value.
19 def encode_text_single_dummy(df, name, target_values):
20     for tv in target_values:
21         l = list(df[name].astype(str))
22         l = [1 if str(x) == str(tv) else 0 for x in l]
23         name2 = "{}-{}".format(name, tv)
24         df[name2] = l
25
26
27 # Encode text values to indexes(i.e. [1],[2],[3] for red,green,blue).
28 def encode_text_index(df, name):
29     le = preprocessing.LabelEncoder()
30     df[name] = le.fit_transform(df[name])
31     return le.classes_
32
33
```

```

34 # Encode a numeric column as zscores
35 def encode_numeric_zscore(df, name, mean=None, sd=None):
36     if mean is None:
37         mean = df[name].mean()
38
39     if sd is None:
40         sd = df[name].std()
41
42     df[name] = (df[name] - mean) / sd
43
44
45 # Convert all missing values in the specified column to the median
46 def missing_median(df, name):
47     med = df[name].median()
48     df[name] = df[name].fillna(med)
49
50
51 # Convert all missing values in the specified column to the default
52 def missing_default(df, name, default_value):
53     df[name] = df[name].fillna(default_value)
54
55
56 # Convert a Pandas dataframe to the x,y inputs that TensorFlow needs
57 def to_xy(df, target):
58     result = []
59     for x in df.columns:
60         if x != target:
61             result.append(x)
62     # find out the type of the target column. Is it really this hard? :(
63     target_type = df[target].dtypes
64     target_type = target_type[0] if hasattr(target_type, '__iter__') else ←
        target_type
65     # Encode to int for classification, float otherwise. TensorFlow likes ←
        32 bits.
66     if target_type in (np.int64, np.int32):
67         # Classification
68         dummies = pd.get_dummies(df[target])
69         return df.as_matrix(result).astype(np.float32), dummies.as_matrix(←
            ().astype(np.float32)
70     else:
71         # Regression
72         return df.as_matrix(result).astype(np.float32), df.as_matrix([←
            target]).astype(np.float32)
73
74 # Nicely formatted time string
75 def hms_string(sec_elapsed):
76     h = int(sec_elapsed / (60 * 60))

```

```

77     m = int((sec_elapsed % (60 * 60)) / 60)
78     s = sec_elapsed % 60
79     return "{:}:{:}>02}:{:}>05.2f}".format(h, m, s)
80
81
82 # Regression chart.
83 def chart_regression(pred,y,sort=True):
84     t = pd.DataFrame({'pred' : pred, 'y' : y.flatten()})
85     if sort:
86         t.sort_values(by=['y'],inplace=True)
87         a = plt.plot(t['y'].tolist(),label='expected')
88         b = plt.plot(t['pred'].tolist(),label='prediction')
89         plt.ylabel('output')
90         plt.legend()
91         plt.show()
92
93 # Remove all rows where the specified column is +/- sd standard deviations
94 def remove_outliers(df, name, sd):
95     drop_rows = df.index[(np.abs(df[name] - df[name].mean()) >= (sd * df[←
        name].std()))]
96     df.drop(drop_rows, axis=0, inplace=True)
97
98
99 # Encode a column to a range between normalized_low and normalized_high.
100 def encode_numeric_range(df, name, normalized_low=-1, normalized_high=1,
101     data_low=None, data_high=None):
102     if data_low is None:
103         data_low = min(df[name])
104         data_high = max(df[name])
105
106     df[name] = ((df[name] - data_low) / (data_high - data_low)) \
107         * (normalized_high - normalized_low) + normalized_low
108
109 # Solution
110
111 def encode_toy_dataset(filename):
112     df = pd.read_csv(filename, na_values=['NA', '?'])
113     encode_numeric_zscore(df, 'length')
114     encode_numeric_zscore(df, 'width')
115     encode_numeric_zscore(df, 'height')
116     encode_text_dummy(df, 'metal')
117     encode_text_dummy(df, 'shape')
118     return df
119
120 # Encode the toy dataset
121 def question1():
122     print()

```

```

123     print("***Question 1***")
124
125     path = "./data/"
126
127     filename_read = os.path.join(path, "toy1.csv")
128     filename_write = os.path.join(path, "submit-jheaton-prog2q1.csv")
129     df = encode_toy_dataset(filename_read) # You just have to implement ↔
        encode_toy_dataset above
130     df.to_csv(filename_write, index=False)
131     print("Wrote {} lines.".format(len(df)))
132
133
134 # Model the toy dataset, no cross validation
135 def question2():
136     print()
137     print("***Question 2***")
138
139 def question3():
140     print()
141     print("***Question 3***")
142
143     # Z-Score encode these using the mean/sd from the dataset (you got ↔
        this in question 2)
144     testDF = pd.DataFrame([
145         {'length':1, 'width':2, 'height': 3},
146         {'length':3, 'width':2, 'height': 5},
147         {'length':4, 'width':1, 'height': 3}
148     ])
149
150
151 def question4():
152     print()
153     print("***Question 4***")
154
155
156 def question5():
157     print()
158     print("***Question 5***")
159
160
161 question1()
162 question2()
163 question3()
164 question4()
165 question5()

```

Listing 2 shows what the output from this assignment would look like. Your numbers might differ from mine slightly. **Every question, except 2, also generates an output CSV file.** For your submission please include your Jupyter notebook and any generated CSV files that the questions specified. Name your output CSV files something such as **submit-jheaton-prog2q1.csv**. Submit a ZIP file that contains your Jupyter notebook and 4 CSV files to Blackboard. This will be 5 files total.

Listing 2: Expected Output

```
1  ***Question 1***
2  Wrote 10001 lines.
3
4  ***Question 2***
5  Epoch 00144: early stopping
6  Final score (RMSE): 75.46247100830078
7
8  ***Question 3***
9  length: (5.5258474152584744, 2.8609014041584113)
10 width: (5.5340465953404658, 2.8598366585224158)
11 height: (5.5337466253374661, 2.8719829476156122)
12      height    length    width
13 0 -0.882205 -1.581907 -1.235659
14 1 -0.185856 -0.882861 -1.235659
15 2 -0.882205 -0.533338 -1.585338
16
17 ***Question 4***
18 Fold #1
19 Epoch 00060: early stopping
20 Fold score (RMSE): 0.21216803789138794
21 Fold #2
22 Epoch 00061: early stopping
23 Fold score (RMSE): 0.14340682327747345
24 Fold #3
25 Epoch 00028: early stopping
26 Fold score (RMSE): 0.3336745500564575
27 Fold #4
28 Epoch 00058: early stopping
29 Fold score (RMSE): 0.2133668214082718
30 Fold #5
31 Epoch 00077: early stopping
32 Fold score (RMSE): 0.1796143352985382
33 Final, out of sample score (RMSE): 0.22570167481899261
34
35 ***Question 5***
36 Fold #1
37 Epoch 00182: early stopping
```

```
38 Fold score: 0.3625
39 Fold #2
40 Epoch 00425: early stopping
41 Fold score: 0.9875
42 Fold #3
43 Epoch 00169: early stopping
44 Fold score: 0.975
45 Fold #4
46 Epoch 00111: early stopping
47 Fold score: 0.8987341772151899
48 Fold #5
49 Epoch 00203: early stopping
50 Fold score: 0.8227848101265823
51 Final, out of sample score: 0.8090452261306532
```

Question 1

Use the dataset found here for this question: [\[click for toy dataset\]](#).

Encode the **toy1.csv** dataset. Generate dummy variables for the shape and metal. Encode height, width and length as z-scores. **Include, but do not encode the weight.** If this encoding is performed in a function, named **encode_toy_dataset**, you will have an easier time reusing the code from question 1 in question 2.

Write the output to a CSV file that you will submit with this assignment. The CSV file will look similar to Listing 3.

Listing 3: Question 2 Output Sample

```
1 height,length,width,weight,metal-bronze,metal-gold,metal-platinum,metal-↵
   silver,metal-tin,shape-box,shape-cylinder,shape-sphere
2 -0.18585564084337075, -0.18381430131795315, -0.1866234261937586, ↵
   729.63,1.0,0.0,0.0,0.0,0.0,0.0,1.0,0.0
3 -0.5340303145851667, -0.18381430131795315, 0.16305509393694917, ↵
   2530.8,0.0,0.0,1.0,0.0,0.0,1.0,0.0,0.0
4 -1.2303796620687586, -1.2323841890415879, -1.5853375067165896, ↵
   58.37,1.0,0.0,0.0,0.0,0.0,0.0,1.0,0.0
5 -0.8822049883269626, -0.8828608931337095, -0.8859804664551741, ↵
   74.15,0.0,0.0,0.0,1.0,0.0,0.0,0.0,1.0
6 ...
```

Question 2

Use the dataset found here for this question: [\[click for toy dataset\]](#).

Use the encoded dataset from question 1 and train a neural network to predict weight. Use 25% of the data as validation and 75% as training, make sure you shuffle the data. Report the RMSE error for the validation set. No CSV file need be generated for this question.

Question 3

Use the dataset found here for this question: [\[click for toy dataset\]](#).

Using the **toy1.csv** dataset calculate and report the mean and standard deviation for height, width and length. Calculate the z-scores for the dataframe given by Listing 4. Make sure that you use the mean and standard deviations you reported for this question. Write the results to a CSV file.

Listing 4: Question 3 Input Data

```
1 testDF = pd.DataFrame([
2     {'length':1, 'width':2, 'height': 3},
3     {'length':3, 'width':2, 'height': 5},
4     {'length':4, 'width':1, 'height': 3}
5 ])
6 ...
```

Your resulting CSV file should look almost exactly like Listing 5.

Listing 5: Question 3 Output Sample

```
1 height,length,width
2 -0.8822049883269626,-1.5819074849494659,-1.2356589865858818
3 -0.18585564084337075,-0.8828608931337095,-1.2356589865858818
4 -0.8822049883269626,-0.5333375972258314,-1.5853375067165896
```

Question 4

Use the dataset found here for this question: [\[click for iris dataset\]](#).

Usually the **iris.csv** dataset is used to classify the species. Not this time! Use the fields species, sepal-l, sepal-w, and petal-l to predict petal-w. Use a 5-fold cross validation and report ONLY out-of-sample predictions to a CSV file. Make sure to shuffle the data. Your generated CSV file should look similar to Listing 6. Encode each of the inputs in a way that makes sense (e.g. dummies, z-scores).

Listing 6: Question 4 Output Sample

```
1 sepal_l,sepal_w,petal_l,petal_w,species-Iris-setosa,species-Iris-↵
   versicolor,species-Iris-virginica,0,0
```

```

2 0.30995914214417364, -0.5903951331558184, 0.5336208818725668, ↵
    1.2,0.0,1.0,0.0,1.2,1.444551944732666
3 -0.1730940663922016, 1.7038864723719687, -1.1658086782311483, ↵
    0.3,1.0,0.0,0.0,0.3,0.\
4 ...

```

Question 5

Use the dataset found here for this question: [click for auto mpg dataset].

Usually the **auto-mpg.csv** dataset is used to regress the mpg. Not this time! Use the fields to predict how many cylinders the car has. Treat this as a classification problem, where there is a class for each number of cylinders. Use a 5-fold cross validation and report ONLY out-of-sample predictions to a CSV file. Make sure to shuffle the data. Your generated CSV file should look similar to Listing 7. Encode each of the inputs in a way that makes sense (e.g. dummies, z-scores). Report the final out of sample accuracy score.

Listing 7: Question 4 Output Sample

```

1 mpg,cylinders,displacement,horsepower,weight,acceleration,year,origin,name↵
    ,ideal,predict
2 -0.7055506566787514, 8, 1.0892327311042995, 0.6722714619460141, ↵
    0.6300768256149949, -1.2938698102195594, 70, -0.7142457922976494,↵
    chevrolet chevelle malibu,8,8
3 -1.0893794720944747, 8, 1.5016242793620063, 1.5879594901955474, ↵
    0.8532590135498572, -1.4751810504376373, 70, -0.7142457922976494,buick↵
    skylark 320,8,8
4 -0.7055506566787514, 8, 1.1947282434492943, 1.19552176380289, ↵
    0.5497784722839334, -1.6564922906557151, 70, -0.7142457922976494,↵
    plymouth satellite,8,8
5 ...

```
