# TripType Classifier on Purchase Information

Lanxin(Lancy) Xu

*Abstract:*

Based on purchase information from Walmart, a classification model is built using Keras to predict TripType(the type of a shopping trip used by Walmart). This process involved necessary data cleansing, model experiments and model tuning. The best model has a multi-class accuracy of 67%.

## Contents

# Objective

As said in Deliverable 1, my goal is to predict the TripType based on purchase information. The original train.csv has 38 different TripType. I have tried to predict 38 classes in my model, but it seems that the large number of classes in response variable has a great influence on accuracy rate. Besides, considering the large size of my data (650k rows), it will be a better choice to predict less classes. Here, I choose to predict TripType '25', '32', '40' and 'others'. This will be explained in Data Section.

# Data

## 1. Data Description

Train.csv: 647,055 rows and 7 columns

Fields:

- TripType - a categorical id representing the type of shopping trip the customer made. This is the ground truth that you are predicting. TripType_999 is an "other" category
- VisitNumber - an id corresponding to a single trip by a single customer
- Weekday - the weekday of the trip
- Upc - the UPC number of the product purchased
- ScanCount - the number of the given item that was purchased. A negative value indicates a product return
- DepartmentDescription - a high-level description of the item's department
- FinelineNumber - a more refined category for each of the products, created by Walmart

Note that Walmart also provides a test.csv file which lacks the TripType for competition. Therefore, it is not used here.

| Field | Type | NumOfClasses | Action |
|---|---|---|---|
| TripType | Category(int) | 38 | Response |
| VisitNumber | id | - | Discard |
| Weekday | Category(string) | 7 | Feature |
| Upc | Category(int) | 93,920 | Discard |
| ScanCount | Numeric | - | Feature |
| DepartmentDescription | Category(string) | 68 | Feature |
| FinelineNumber | Category(int) | 5195 | Feature |

**Form 1.1-Field Description**

## 2. Data Issue

The data suffers from several problems which need to be solve before moving on to model building step.

## a) Multi-class Response Variable 'TripType'

As mentioned before, the TripType field has 38 unique classes which, if transformed into dummy variables, will result in 38 columns as response variable. Such a large number of classes in y will have a negative effect on prediction accuracy. Therefore, I need to reduce it.
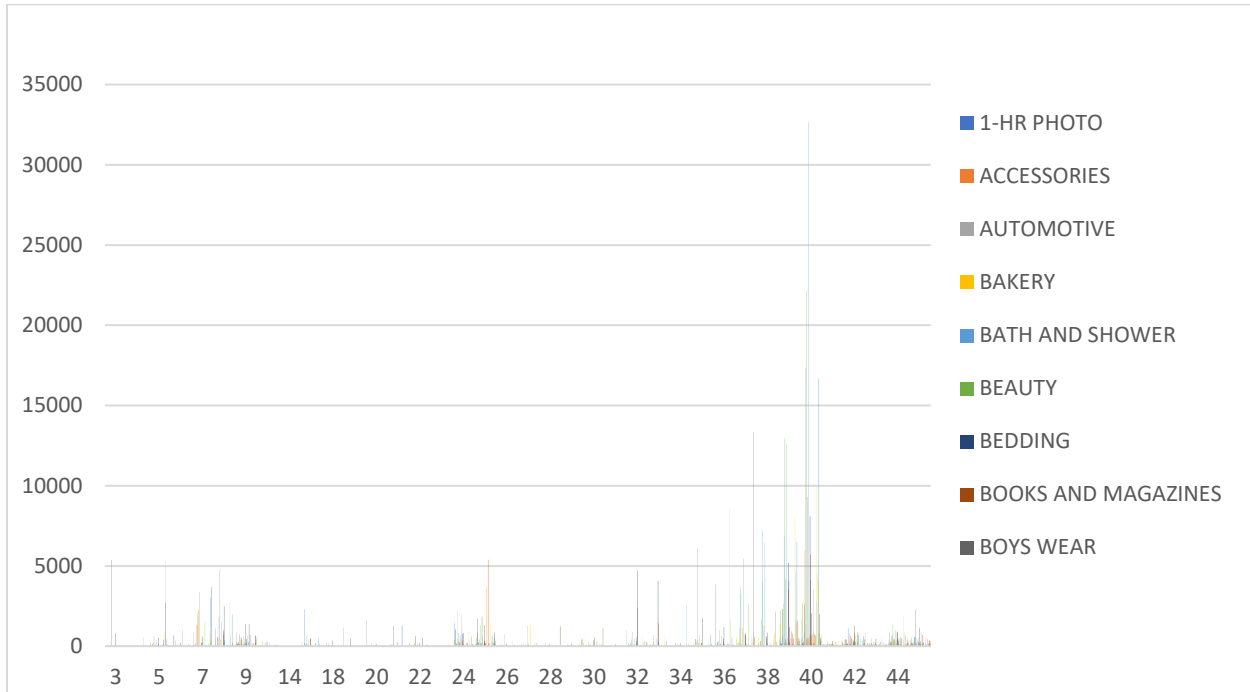


**Figure 2.1-DepartmentDescription by TripType**

Figure 2.1 shows a plot between DepartmentDescription and TripType. It is easy to find that some TripType has an obvious relationship with particular Department. For example, Type25 has a highest number of shopping trips in 'MENS_WEAR', this triptype maybe mainly related to male shoppers. After investigating some patterns and considered the size of each TripType, I choose Type '25', '32', '40' as my main research objects and set all other TripType to 999 (means 'others'). Therefore, the response variable only has 4 classes.

## b) Missing/Null Value

In train.csv file, there are missing values appear in column 'Upc' and 'FinelineNumber'. Besides, in 'DepartmentDescription', there are some records shown as NULL which is meaningless. Therefore, in Excel, I checked these three columns and found that there are only 4,129 out of 647,054 records that have one or more problems. So, to simplify the data cleaning process, I just delete those records.

**Figure 2.2-Check for problematic records**

## c) Feature Selection

As shown in Form 1.1, not all fields are used in the model. Some are discarded due to multiple reasons:

VisitNumber- an id number. No relationship with TripType.

Upc- an id for product purchased. If used in the model, will results in 93920 dummy variables. The large number of predictors will lead to serious overfitting problem. However, to take full use of the data, I did some experiments on UPC. (See Experiments Section)

## 3. Data Preprocessing

### a) Data Encoding and Dummy Variables

All the category variables need to be treated differently with numeric variables in model. Since weekday and DepartmentDescription are string, they need to be transformed into numbers. I used sklearn.preprocessing.LabelEncoder to finish this work.

```
def Encode_X(X,area):
    for i in area:
        labelencoder_X = LabelEncoder()
        X[:, i] = labelencoder_X.fit_transform(X[:, i])
    return X
```

Then, I used sklearn.preprocessing.OneHotEncoder to create n dummy variables for n classes for each variable.

```
def create_dummy(X,area):
    onehotencoder = OneHotEncoder(categorical_features = area)
    X = onehotencoder.fit_transform(X)
    return X
```

## b) Scaling

To prevent that some variables with larger values will have larger weights in model, I scale all my features using preprocessing.MaxAbsScaler. The choice of scaler was limited to me because my feature matrix is a sparse matrix (can't be turned into dense due to memory issue). MinMax Scaler is not suitable for sparse matrix, so I choose its similar scaler MaxAbs to scale my data.

## c) Data slicing

To test my model, I slice my data into train and test setsDue to the large size of data, it requires much time to tune on test size. So I just set the test size as 0.3 to simplify the process.

# Model

## 1. Overview

| Data Input | NB_units+ activation | Dropout (Dropout_rate) | NB_units+ activation | Dropout (Dropout_rate) | 4 units+ Softmax |
|---|---|---|---|---|---|

## 2. Model Tuning

To find the best parameters for my model, I used GridSearchCV to tune my model:

```
def tuning(model,param_grid,x_train,y_train):
    grid = GridSearchCV(estimator=model, param_grid=param_grid)
    grid_result = grid.fit(x_train, y_train,verbose=1)
    return grid_result
```

I made a param_grid with the range of each parameter I want to tune.

```
neurons = [50,100,150]
dropout_rate = [0.0, 0.1, 0.2]
activation = ['relu','sigmoid']
optimizer = ['SGD', 'RMSprop','Adam']
batch_size = [10, 50, 100]
epochs = [10, 20, 50]
param_grid = dict(batch_size=batch_size, epochs=epochs,
                optimizer=optimizer,activation=activation,
                dropout_rate=dropout_rate,neurons=neurons)
```

## 3. Result and Evaluation

### a) Loss and Metrics

Since my response variable is a multi-class variable, I choose categorical_crossentropy as loss and set metrics equals to accuracy

### b) Result

*i. Best model*

| Data Input | → | 150 units+ Sigmoid | → | Dropout 0.1 | → | 150 units+ Sigmoid | → | Dropout 0.1 | → | 4 units+ Softmax |

Batch_size=100,

Epochs=50

*ii. Evaluation*

Time required for optimization: 0:25:43.586648

Test Accuracy = 0.673721212374

Here is the summarized history of accuracy and loss:



**Figure 3.3.1- Summarized history of accuracy**



**Figure 3.3.2- Summarized history of loss**

A model with accuracy of 0.67 normally is not a good one. This might result from the big number of dummy variables, small number of predictors, etc. Also, the large size of data is a great challenge in this assignment.

# Experiments

## 1. UPC

Always want to use as much data as I can, I did some experiments on Upc. Although they all failed in the end, it is good to show my work and thoughts there.

In my data, there is a column 'Upc' which shows the UPC number of the product purchased. I did some research and found that they are not in the standard UPC-A format. In short, UPCs are Universal Product Code that consists of the company code, a product number and a checksum.



**Figure 4.1-Anatomy of a UPC Symbol**

I calculate the number of digits and found the results below:

| Row Labels | Count of Number of Digits |
|---|---|
| 0 | 4129 |
| 3 | 7 |
| 4 | 29745 |
| 5 | 372 |
| 7 | 1 |
| 8 | 412 |
| 9 | 2166 |
| 10 | 433341 |
| 11 | 168418 |
| 12 | 8463 |
| **Grand Total** | **647054** |

**Figure 4.2-Check the number of digits in 'Upc'**

I read some posts from Kaggle discussion board and it seemed that most of the participants agreed that most of the UPCs that are less than 12 digits lack the check digit at the end. I don't need to calculate the check digit since my goal is to get the Company_code and Product_code so that I can add variable in my classification.

Therefore, to simplify my analysis, I assume that all Company_code has 6 digits and the rest are Product_code. There are some records with 3,4 and 5 digits, which, to me are some in-house products. So I set their Company_code as 0 and their Upc as their Product_code. For missing values(number of digit=0), I set both variables as -999.
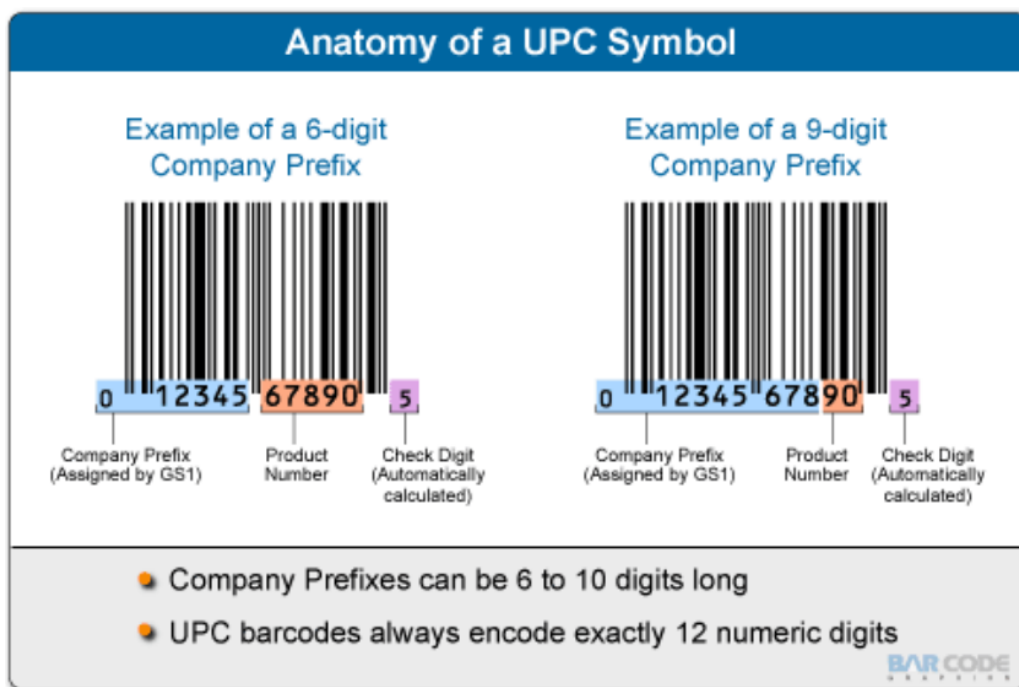
By doing all of the process above, I changed the 'Upc' variable into 'Company_code' and 'Product_code'.

However, again, due to the large number of classes in Company_code and Product_code, they were not used in the model. But it is still possible that if given more information (e.g. the name of Company), I may be able to combine similar classes to reduce the dimentionality.

```
Data Input → 150 units+ Sigmoid → Dropout 0.1 → 150 units+ Sigmoid → Dropout 0.1 → 4 units+ Softmax
```

## 2. Embedding layers

When I met with the problem of too many dummy variables, I searched online and found that Embedding Layers which provided by Keras could give me some help.

According to Keras, embedding layers can be used to turn positive integers (indexes) into dense vectors of fixed size. It is used a lot in text mining where there are many words need to be transformed into numeric values in model. Apparently, it is impractical to use dummy variables.

Therefore, I try to use embedding layers to deal with the category features:

Category Input
(Input Layer)

↓

Embedding Layer

↓

Numeric Input
(Input Layer)          LSTM

↘          ↙

Merge

↓

Dense

↓

Dropout

↓

Dense

↓

Dropout

↓

Output

However, the accuracy rate doesn't have significant improvement. Besides, since I don't fully understand the mechanism behind embedding layers, I just decide not to use it as my final model.

## Discussion

Dealing with multi-class features like FinelineNumber(about 6000 classes) or Upc(90k) is really the most difficult part. I tried to turn them all into dummy variables and resulted in a matrix of 50k columns. Using this data to train my model, I got an extremely low accuracy rate of 0.25. This poor performance forced me to discard Upc but also made me think whether there is a way to keep my multi-classes features in my model. Because whenever we train the model, we would like to feed it with more information and take full advantage of the data that we have. As I mentioned before, one way is to embed them to reduce dimensionality. Inspired by text mining, I also came up with an idea of using the techniques in

that field to deal with the category variables. After all, the different categories (whether string or numeric) are essentially words. Just as people use word2vec in text mining to reduce dimensionality, we can apply the same principle in category. I am glad to find that many people have the same idea with me. There is already some research and experiments (e.g. a similar method called cat2vec) in this field and will be more improvement in the future. Discussion and further work.

# Code

```python
1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Aug 30 14:20:41 2016
4
5  @author: Lancy(Lanxin) Xu
6  """
7
8  """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
9  Import Libraries Section
10 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
11
12 #from __future__ import print_function
13 from sklearn import preprocessing
14 from sklearn.model_selection import train_test_split
15 from sklearn.model_selection import GridSearchCV
16 import random
17 import matplotlib.pyplot as plt
18 import datetime
19 import numpy as np
20 import pandas as pd
21
22 import keras
23 from keras.wrappers.scikit_learn import KerasClassifier
24 from keras.models import Sequential
25 from keras.layers import Dense, Dropout, Activation, Input, Embedding, Merge, LSTM
26 from keras.models import Model
27 from keras.optimizers import RMSprop,SGD, Adam
28 from keras.utils import np_utils
29 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
30 from scipy.sparse import csr_matrix
31
32 np.random.seed(2018)
33 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
34 Define Function
35 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
36 def load_data(URL):
37     data=pd.read_csv(URL,
38                 delimiter=",", header =0,
39                 names=['TripType','Weekday','DepartmentDescription','FinelineNumber','ScanCount','Check_problem'])
40     return data

41
42 def preprocess(data):
43     check_problem=data[data.Check_problem==0] #Delete all rows with missing values or with Null values
44     drop_check=check_problem.drop('Check_problem',1) #Drop the Check-problem column since all problematic rows were deleted
45     return drop_check
46
47 def get_NB_CLASSES(data):
48     NB_CLASSES=len(pd.Series.unique(data.TripType))
49     #The number of output classes should be the number of unique classes in response variable
50     return NB_CLASSES
51
52 def Encode_X(X,area):
53     for i in area:
54         labelencoder_X = LabelEncoder()
55         X[:, i] = labelencoder_X.fit_transform(X[:, i])
56     return X
57
58 def Encode_Y(Y):
59     labelencoder_Y = LabelEncoder()
60     Y= labelencoder_Y.fit_transform(Y)
61     Y=np_utils.to_categorical(Y)
62     return Y
63
64 def create_dummy(X,area):
65     onehotencoder = OneHotEncoder(categorical_features = area)
66     X = onehotencoder.fit_transform(X)
67     return X
68
69 def MaxAbs(X):
70     X_MaxAbs = preprocessing.MaxAbsScaler()
71     X = X_MaxAbs.fit_transform(X)
72     return X
```

```python
73
74 def create_model(optimizer='RMSprop',activation='relu',dropout_rate=0.1,neurons=50):
75     model = Sequential()
76     model.add(Dense(neurons,input_shape=(x_train.shape[1],)))
77     model.add(Activation(activation))
78     model.add(Dropout(dropout_rate))
79     model.add(Dense(neurons))
80     model.add(Activation(activation))
81     model.add(Dropout(dropout_rate))
82     model.add(Dense(NB_CLASSES))
83     model.add(Activation('softmax'))
84     model.compile(loss='categorical_crossentropy', optimizer=optimizer, metrics=['accuracy'])
85     return model
86
87 def tuning(model,param_grid,x_train,y_train):
88     grid = GridSearchCV(estimator=model, param_grid=param_grid)
89     grid_result = grid.fit(x_train, y_train,verbose=1)
90     return grid_result
91
92 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
93 Set Parameters
94 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
95 VALIDATION_SPLIT=0.3
96 TEST_SIZE=0.3
97
98 neurons = [50,100,150]
99 dropout_rate = [0.0, 0.1, 0.2]
100 activation = ['relu','sigmoid']
101 optimizer = ['SGD', 'RMSprop','Adam']
102 batch_size = [10, 50, 100]
103 epochs = [10, 20, 50]
104 param_grid = dict(batch_size=batch_size, epochs=epochs,
105                   optimizer=optimizer,activation=activation,
106                   dropout_rate=dropout_rate,neurons=neurons)
107
108
109 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
110 Load Data Section
111 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
112 #Import trainset
113 data=load_data("D:/myfile/Spring/Artifitial Intelligence/Final/train.csv")
114 #Preprocess
115 data=preprocess(data)
116
117 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
118 Pretreat Data Section
119 """"""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
120 x=data[data.columns[1:]].values
121 y=data.TripType.values
122 NB_CLASSES=len(pd.Series.unique(data.TripType))
123 #Encode x
124 area=[0,1,2]
125 x=Encode_X(x,area)
126 #Encode y
127 y=Encode_Y(y)
128
129 #Create Dummy variables
130 x=create_dummy(x,area)
131 x=csr_matrix(x)
132
133 #Scale
134 #Using MaxAbs because my x is in a sparse matrix and can not be turned into dense matrix due to memory issue.
135 #Therefore, I can not use MinMax scaler.
136 x=MaxAbs(x)
137
138 #Data slicing
139 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = TEST_SIZE)
140
```

```python
141 """""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
142 Define Model Section
143 """""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
144 #Build Model
145 model=KerasClassifier(build_fn=create_model, verbose=1)
146
147 #Tune model
148 grid = GridSearchCV(estimator=model, param_grid=param_grid)
149 grid_result = grid.fit(x_train, y_train,verbose=1)
150
151 print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
152
153 """""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
154 Show output Section
155 """""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""""
156 start_time = datetime.datetime.now()
157
158 #Below shows the best_params
159 neurons = 150
160 dropout_rate = 0.1
161 activation = 'sigmoid'
162 optimizer = 'Rmsprop'
163 batch_size = 100
164 epochs = 50
165
166 #Build model
167 model = Sequential()
168 model.add(Dense(neurons,input_shape=(x_train.shape[1],)))
169 model.add(Activation(activation))
170 model.add(Dropout(dropout_rate))
171 model.add(Dense(neurons))
172 model.add(Activation(activation))
173 model.add(Dropout(dropout_rate))
174 model.add(Dense(NB_CLASSES))
175 model.add(Activation('softmax'))
176 model.compile(loss='categorical_crossentropy', optimizer=optimizer, metrics=['accuracy'])
177
178 history = model.fit(x_train, y_train,
179                     epochs=epochs, batch_size=batch_size,
180                     verbose=1, validation_split=0.3)
181
182 print('Testing...')
183
184 #Evaluate Model
185 score = model.evaluate(x_test, y_test, verbose=1)
186 print('Test accuracy:', score[1])
187
188 # summarize history for accuracy
189 plt.plot(history.history['acc'])
190 plt.plot(history.history['val_acc'])
191 plt.title('model accuracy')
192 plt.ylabel('accuracy')
193 plt.xlabel('epoch')
194 plt.legend(['train', 'test'], loc='upper left')
195 plt.show()
196
197 # summarize history for loss
198 plt.plot(history.history['loss'])
199 plt.plot(history.history['val_loss'])
200 plt.title('model loss')
201 plt.ylabel('loss')
202 plt.xlabel('epoch')
203 plt.legend(['train', 'test'], loc='upper left')
204 plt.show()
205
206 stop_time = datetime.datetime.now()
207 print ("Time required for optimization:",stop_time - start_time)
```