

Digit Recognition Random Forest Starter

December 31, 2019

```
[1]: import numpy as np

# keras import for the dataset
from keras.datasets import mnist
```

Using TensorFlow backend.

```
[2]: (X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```
[3]: # each training and test element is a 28 x 28 pixel grayvalue image
print(X_train[0].shape)
print(X_train[0])
```

(28, 28)

```
[[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
   0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0  0  3  18  18  18 126 136
 175 26 166 255 247 127  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  30 36 94 154 170 253 253 253 253
225 172 253 242 195 64  0  0  0  0]
 [ 0  0  0  0  0  0  0  49 238 253 253 253 253 253 253 253 253
93 82 82 56 39  0  0  0  0  0]
 [ 0  0  0  0  0  0  0 18 219 253 253 253 253 253 198 182 247 241
  0  0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0 80 156 107 253 253 205 11  0 43 154
  0  0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0 14  1 154 253 90  0  0  0  0
  0  0  0  0  0  0  0  0  0  0]
 [ 0  0  0  0  0  0  0  0  0  0  0 139 253 190  2  0  0  0
  0  0  0  0  0  0  0  0  0  0]
```

```

[ 0  0  0  0  0  0  0  0  0  0  0  11 190 253  70  0  0  0
  0  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  35 241 225 160 108  1
  0  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  81 240 253 253 119
 25  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  45 186 253 253
150 27  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  16  93 252
253 187  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0 249
253 249 64  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  46 130 183 253
253 207  2  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  39 148 229 253 253 253
250 182  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  24 114 221 253 253 253 253 201
 78  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0 23  66 213 253 253 253 253 198  81  2
  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0 18 171 219 253 253 253 253 195  80  9  0  0
  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  55 172 226 253 253 253 253 244 133  11  0  0  0  0
  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0 136 253 253 253 212 135 132  16  0  0  0  0  0  0
  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
  0  0  0  0  0  0  0  0  0  0]
[ 0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
  0  0  0  0  0  0  0  0  0  0]

```

```

[10]: # the corresponding label is the "real" digit
print(np.unique(y_train, return_counts=True))
print(y_train[0])

```

```

(array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9], dtype=uint8), array([5923, 6742, 5958,
6131, 5842, 5421, 5918, 6265, 5851, 5949]))
5

```

```

[5]: # imports for plotting
import matplotlib
matplotlib.use('agg')
import matplotlib.pyplot as plt

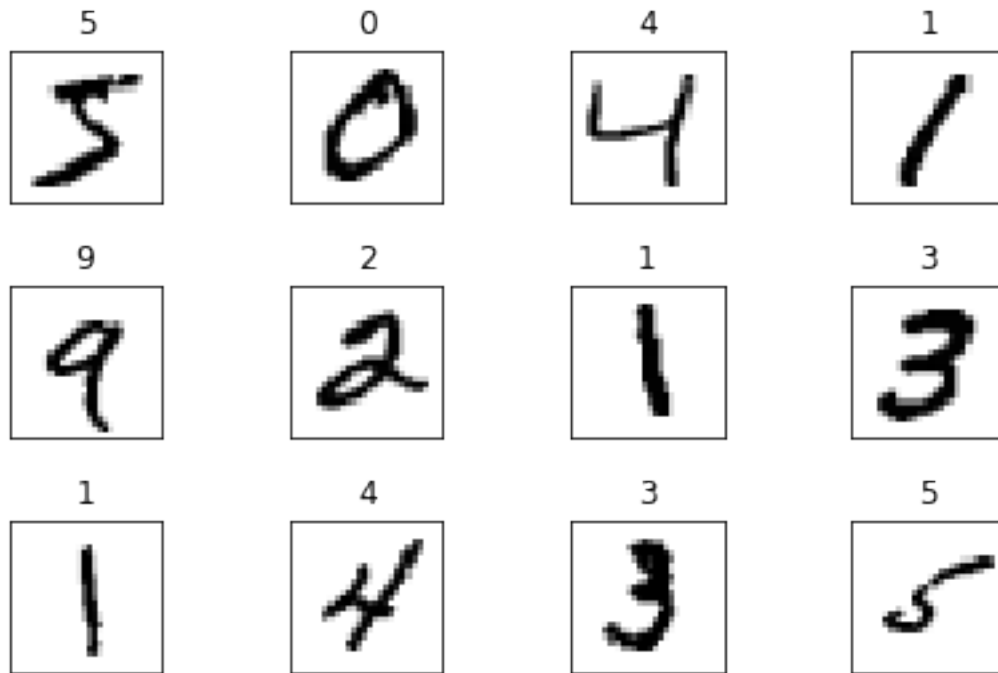
%matplotlib inline

```

```

fig = plt.figure()
for i in range(12):
    plt.subplot(3,4,i+1) # i+1 is position of subplot in 3 x 4 table
    plt.tight_layout() # "tight" layout
    # show bitmap, interpret 0 as white and 255 as black (grayvalues)
    plt.imshow(X_train[i], cmap=plt.cm.gray_r)
    plt.title(y_train[i]) # real value as title
    plt.xticks([]) # no ticks on x axis
    plt.yticks([]) # not ticks on y axis

```



```

[6]: # let's print the shape before we reshape and normalize
print("X_train shape", X_train.shape)
print("y_train shape", y_train.shape)
print("X_test shape", X_test.shape)
print("y_test shape", y_test.shape)

# building the input vector from the 28x28 pixels = linearize the image to get
→ a 784 (= 28x28) vector
X_train = X_train.reshape(60000, 784)
X_test = X_test.reshape(10000, 784)

# normalizing the data to help with the training
# normalized data leads to better models
X_train = X_train.astype('float32')

```

```
X_test = X_test.astype('float32')
X_train /= 255
X_test /= 255

# print the final input shape ready for training
print("Train matrix shape", X_train.shape)
print("Test matrix shape", X_test.shape)
```

```
X_train shape (60000, 28, 28)
y_train shape (60000,)
X_test shape (10000, 28, 28)
y_test shape (10000,)
Train matrix shape (60000, 784)
Test matrix shape (10000, 784)
```

```
[7]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(n_estimators=300)
model.fit(X_train, y_train)
```

```
[7]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                             max_depth=None, max_features='auto', max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=1, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, n_estimators=300,
                             n_jobs=None, oob_score=False, random_state=None,
                             verbose=0, warm_start=False)
```

```
[8]: y_test2 = model.predict(X_test)
```

```
[9]: from sklearn.metrics import accuracy_score
accuracy_score(y_test, y_test2)
```

```
[9]: 0.9711
```

```
[ ]:
```