

[Click to get the 20-book Super Bundle! \(Save \\$250\)](#)

Navigation



Machine Learning Mastery

Making Developers Awesome at Machine Learning

[Click to Take the FREE Computer Vision Crash-Course](#)

Search...



How to Develop a CNN for MNIST Handwritten Digit Classification

by Jason Brownlee on May 8, 2019 in Deep Learning for Computer Vision

[Tweet](#)[Share](#)[Share](#)

Last Updated on August 24, 2020

How to Develop a Convolutional Neural Network From Scratch for MNIST Handwritten Digit Classification.

The MNIST handwritten digit classification problem is a standard dataset used in computer vision and deep learning.

Although the dataset is effectively solved, it can be used as the basis for learning and practicing how to develop, evaluate, and use convolutional deep learning neural networks for image classification from scratch. This includes how to develop a robust test harness for estimating the performance of the model, how to explore improvements to the model, and how to save the model and later load it to make predictions on new data.

In this tutorial, you will discover how to develop a convolutional neural network for handwritten digit classification from scratch.

After completing this tutorial, you will know:

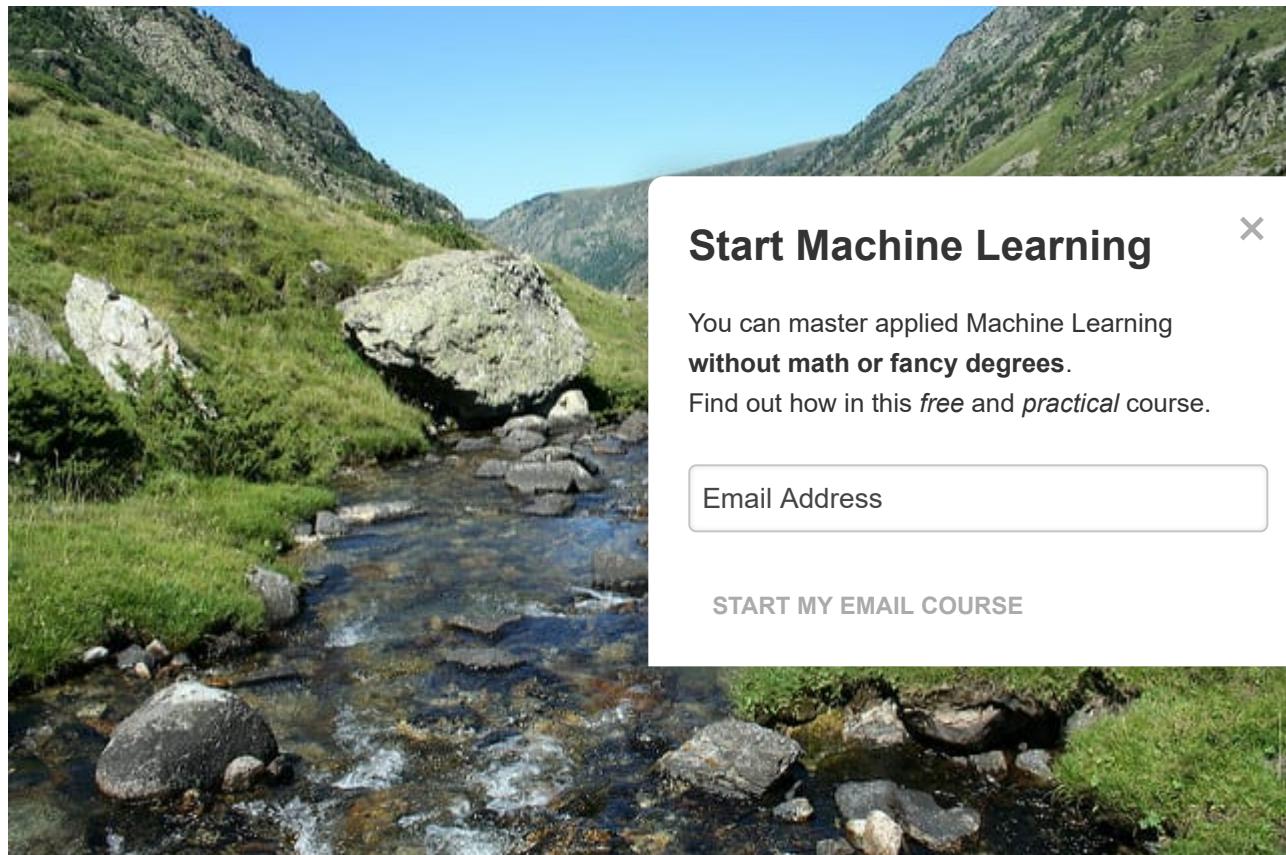
- How to develop a test harness to develop a robust evaluation of a model and establish a baseline of performance for a classification task.
- How to explore extensions to a baseline model to improve learning and model capacity.
- How to develop a finalized model, evaluate the performance of the final model, and use it to make predictions on new images.

[Start Machine Learning](#)

Kick-start your project with my new book [Deep Learning for Computer Vision](#), including *step-by-step tutorials* and the *Python source code* files for all examples.

Let's get started.

- **Updated Dec/2019:** Updated examples for TensorFlow 2.0 and Keras 2.3.
- **Updated Jan/2020:** Fixed a bug where models were defined outside the cross-validation loop.



Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

How to Develop a Convolutional Neural Network From Scratch for MNIST Handwritten Digit Classification
Photo by [Richard Allaway](#), some rights reserved.

Tutorial Overview

This tutorial is divided into five parts; they are:

1. MNIST Handwritten Digit Classification Dataset
2. Model Evaluation Methodology
3. How to Develop a Baseline Model
4. How to Develop an Improved Model
5. How to Finalize the Model and Make Predictions

Want Results with Deep Learning for Computer Vision?

[Take my free 7-day email crash](#)

[Start Machine Learning](#)

Click to sign-up and also get a free PDF Ebook version of the course.

[Download Your FREE Mini-Course](#)

Development Environment

This tutorial assumes that you are using standalone Keras running on top of TensorFlow with Python 3. If you need help setting up your development environment see this tutorial:

- [How to Setup Your Python Environment for Machine Learning with Anaconda](#)

MNIST Handwritten Digit Classification

The [MNIST dataset](#) is an acronym that stands for the [Modified National Institute of Standards and Technology](#) dataset.

It is a dataset of 60,000 small square 28×28 pixel grayscale images of handwritten digits from 0 to 9.

The task is to classify a given image of a handwritten digit into one of ten values from 0 to 9, inclusively.

It is a widely used and deeply understood dataset and, for the most part, is “*solved*.” Top-performing models are deep learning convolutional neural networks that achieve a classification accuracy of above 99%, with an error rate between 0.4 %and 0.2% on the hold out test dataset.

The example below loads the MNIST dataset using the Keras API and creates a plot of the first nine images in the training dataset.

```

1 # example of loading the mnist dataset
2 from keras.datasets import mnist
3 from matplotlib import pyplot
4 # load dataset
5 (trainX, trainy), (testX, testy) = mnist.load_data()
6 # summarize loaded dataset
7 print('Train: X=%s, y=%s' % (trainX.shape, trainy.shape))
8 print('Test: X=%s, y=%s' % (testX.shape, testy.shape))
9 # plot first few images
10 for i in range(9):
11     # define subplot
12     pyplot.subplot(330 + 1 + i)
13     # plot raw pixel data
14     pyplot.imshow(trainX[i], cmap=pyplot.get_cmap('gray'))
15 # show the figure
16 pyplot.show()
```

Running the example loads the MNIST train and test datasets and plots the first nine images in the training dataset.

Start Machine Learning

You can master applied Machine Learning [without math or fancy degrees](#).

Find out how in this *free* and *practical* course.

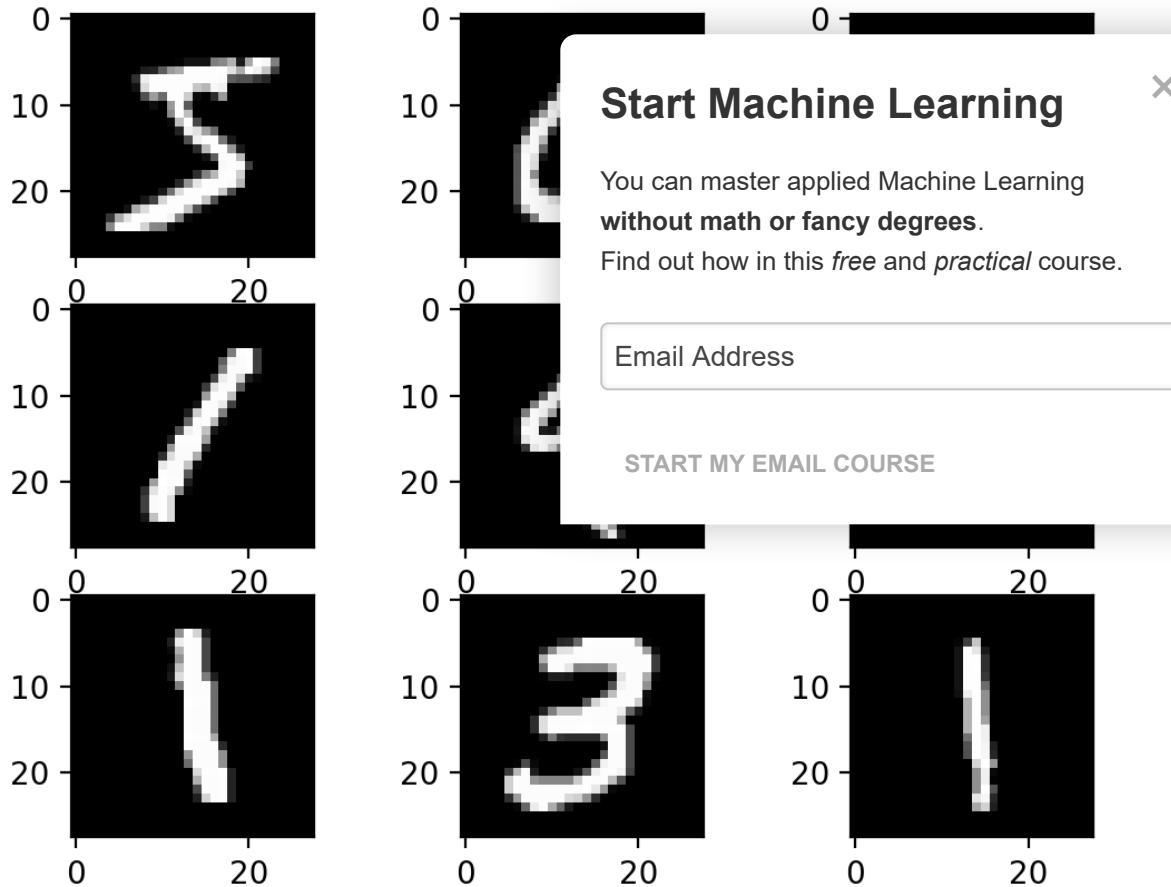
Email Address

[START MY EMAIL COURSE](#)

We can see that there are 60,000 examples in the training dataset and 10,000 in the test dataset and that images are indeed square with 28×28 pixels.

```
1 Train: X=(60000, 28, 28), y=(60000,)Test: X=(10000, 28, 28), y=(10000,)
```

A plot of the first nine images in the dataset is also created showing the natural handwritten nature of the images to be classified.



Plot of a Subset of Images From the MNIST Dataset

Model Evaluation Methodology

Although the MNIST dataset is effectively solved, it can be a useful starting point for developing and practicing a methodology for solving image classification tasks using convolutional neural networks.

Instead of reviewing the literature on well-performing models on the dataset, we can develop a new model from scratch.

The dataset already has a well-defined train and test dataset that we can use

[Start Machine Learning](#)

In order to estimate the performance of a model for a given training run, we can further split the training set into a train and validation dataset. Performance on the train and validation dataset over each run can then be plotted to provide learning curves and insight into how well a model is learning the problem.

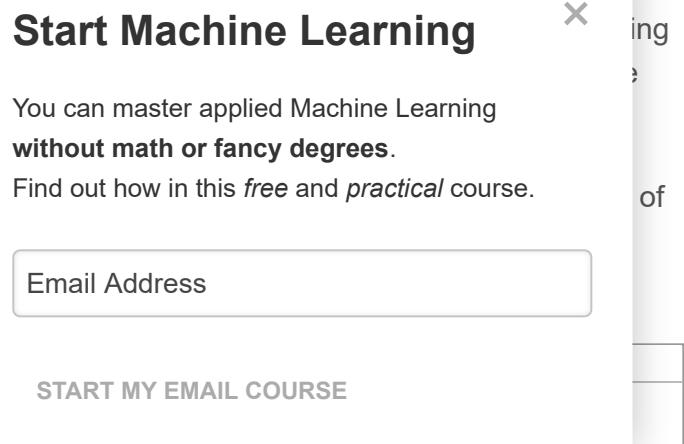
The Keras API supports this by specifying the “*validation_data*” argument to the *model.fit()* function when training the model, that will, in turn, return an object that describes model performance for the chosen loss and metrics on each training epoch.

```
1 # record model performance on a validation dataset during training
2 history = model.fit(..., validation_data=(valX, valY))
```

In order to estimate the performance of a model on the problem in general, we can use **k-fold cross-validation**, perhaps five-fold cross-validation. This will respect to differences in the training and test datasets algorithm. The performance of a model can be taken a standard deviation, that could be used to estimate a confidence interval.

We can use the **KFold** class from the scikit-learn API to perform k-fold cross-validation on a given neural network model. There are many ways to implement this approach where the *KFold* class is only used to specify

```
1 # example of k-fold cv for a neural net
2 data = ...
3 # prepare cross validation
4 kfold = KFold(5, shuffle=True, random_state=1)
5 # enumerate splits
6 for train_ix, test_ix in kfold.split(data):
7     model = ...
8     ...
```



We will hold back the actual test dataset and use it as an evaluation of our final model.

How to Develop a Baseline Model

The first step is to develop a baseline model.

This is critical as it both involves developing the infrastructure for the test harness so that any model we design can be evaluated on the dataset, and it establishes a baseline in model performance on the problem, by which all improvements can be compared.

The design of the test harness is modular, and we can develop a separate function for each piece. This allows a given aspect of the test harness to be modified or inter-changed, if we desire, separately from the rest.

We can develop this test harness with five key elements. They are the loading of the dataset, the preparation of the dataset, the definition of the model, the evaluation of the model, and the presentation of results.

[Start Machine Learning](#)

Load Dataset

We know some things about the dataset.

For example, we know that the images are all pre-aligned (e.g. each image only contains a hand-drawn digit), that the images all have the same square size of 28×28 pixels, and that the images are grayscale.

Therefore, we can load the images and reshape the data arrays to have a single color channel.

```
1 # load dataset
2 (trainX, trainY), (testX, testY) = mnist.load_data()
3 # reshape dataset to have a single channel
4 trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
5 testX = testX.reshape((testX.shape[0], 28, 28, 1))
```

We also know that there are 10 classes and that class

We can, therefore, use a one hot encoding for the clas

into a 10 element binary vector with a 1 for the index c

We can achieve this with the `to_categorical()` utility fun

```
1 # one hot encode target values
2 trainY = to_categorical(trainY)
3 testY = to_categorical(testY)
```

The `load_dataset()` function implements these behaviors and can be used to load the dataset.

```
1 # load train and test dataset
2 def load_dataset():
3     # load dataset
4     (trainX, trainY), (testX, testY) = mnist.load_data()
5     # reshape dataset to have a single channel
6     trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
7     testX = testX.reshape((testX.shape[0], 28, 28, 1))
8     # one hot encode target values
9     trainY = to_categorical(trainY)
10    testY = to_categorical(testY)
11    return trainX, trainY, testX, testY
```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

Prepare Pixel Data

We know that the pixel values for each image in the dataset are unsigned integers in the range between black and white, or 0 and 255.

We do not know the best way to scale the pixel values for modeling, but we know that some scaling will be required.

A good starting point is to normalize the pixel values of grayscale images, e.g. rescale them to the range [0,1]. This involves first converting the data type from unsigned integers to floats, then dividing the pixel values by the maximum value.

[Start Machine Learning](#)

```
1 # convert from integers to floats
2 train_norm = train.astype('float32')
3 test_norm = test.astype('float32')
4 # normalize to range 0-1
5 train_norm = train_norm / 255.0
6 test_norm = test_norm / 255.0
```

The `prep_pixels()` function below implements these behaviors and is provided with the pixel values for both the train and test datasets that will need to be scaled.

```
1 # scale pixels
2 def prep_pixels(train, test):
3     # convert from integers to floats
4     train_norm = train.astype('float32')
5     test_norm = test.astype('float32')
6     # normalize to range 0-1
7     train_norm = train_norm / 255.0
8     test_norm = test_norm / 255.0
9     # return normalized images
10    return train_norm, test_norm
```

This function must be called to prepare the pixel values.

Define Model

Next, we need to define a baseline convolutional neural network.

The model has two main aspects: the feature extraction front-end comprised of convolutional and pooling layers, and the classifier backend that will make a prediction.

For the convolutional front-end, we can start with a single convolutional layer with a small filter size (3,3) and a modest number of filters (32) followed by a max pooling layer. The filter maps can then be flattened to provide features to the classifier.

Given that the problem is a multi-class classification task, we know that we will require an output layer with 10 nodes in order to predict the probability distribution of an image belonging to each of the 10 classes. This will also require the use of a softmax activation function. Between the feature extractor and the output layer, we can add a dense layer to interpret the features, in this case with 100 nodes.

All layers will use the [ReLU activation function](#) and the He weight initialization scheme, both best practices.

We will use a conservative configuration for the stochastic gradient descent optimizer with a [learning rate](#) of 0.01 and a momentum of 0.9. The [categorical cross-entropy](#) loss function will be optimized, suitable for multi-class classification, and we will monitor the classification accuracy metric, which is appropriate given we have the same number of examples in each of the 10 classes.

The `define_model()` function below will define and return this model.

```
1 # define cnn model
```

Start Machine Learning

```

1 def define_model():
2     model = Sequential()
3     model.add(Conv2D(32, (3, 3), activation='relu', kernel_initializer='he_uniform', input_shape=(28, 28, 1)))
4     model.add(MaxPooling2D((2, 2)))
5     model.add(Flatten())
6     model.add(Dense(100, activation='relu', kernel_initializer='he_uniform'))
7     model.add(Dense(10, activation='softmax'))
8
9     # compile model
10    opt = SGD(lr=0.01, momentum=0.9)
11    model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
12
13    return model

```

Evaluate Model

After the model is defined, we need to evaluate it.

The model will be evaluated using **five-fold cross-validation**. This provides a baseline for both repeated evaluation and to not be so sensitive to the specific training and test sets. The validation set will be 20% of the training dataset, or about 12,000 examples. This is a good size for this problem.

The training dataset is shuffled prior to being split, and the same splits will be used for any model we evaluate. This means that each model we evaluate will have the same train and test sets. This allows for a fair apples comparison between models.

We will train the baseline model for a modest 10 training epochs. We will also keep the validation set small. The test set for each fold will be used to evaluate the model. This provides a baseline for the model's performance. It also allows us to track the model's performance over time, so that we can later create learning curves, and at the end of the run, so that we can estimate the performance of the model. As such, we will keep track of the resulting history from each run, as well as the classification accuracy of the fold.

The `evaluate_model()` function below implements these behaviors, taking the training dataset as arguments and returning a list of accuracy scores and training histories that can be later summarized.

```

1 # evaluate a model using k-fold cross-validation
2 def evaluate_model(dataX, dataY, n_folds=5):
3     scores, histories = list(), list()
4     # prepare cross validation
5     kfold = KFold(n_folds, shuffle=True, random_state=1)
6     # enumerate splits
7     for train_ix, test_ix in kfold.split(dataX):
8         # define model
9         model = define_model()
10        # select rows for train and test
11        trainX, trainY, testX, testY = dataX[train_ix], dataY[train_ix], dataX[test_ix], dataY[test_ix]
12        # fit model
13        history = model.fit(trainX, trainY, epochs=10, batch_size=32, validation_data=(testX, testY))
14        # evaluate model
15        _, acc = model.evaluate(testX, testY, verbose=0)
16        print('> %.3f' % (acc * 100.0))
17        # store scores
18        scores.append(acc)
19        histories.append(history)
20
21    return scores, histories

```

[Start Machine Learning](#)

Present Results

Once the model has been evaluated, we can present the results.

There are two key aspects to present: the diagnostics of the learning behavior of the model during training and the estimation of the model performance. These can be implemented using separate functions.

First, the diagnostics involve creating a line plot showing model performance on the train and test set during each fold of the k-fold cross-validation. These plots are valuable for getting an idea of whether a model is overfitting, underfitting, or has a good fit for the dataset.

We will create a single figure with two subplots, one for model performance on the training dataset and orange for the test dataset. The `summarize_diagnostics()` function below summarizes training histories.

```

1 # plot diagnostic learning curves
2 def summarize_diagnostics(histories):
3     for i in range(len(histories)):
4         # plot loss
5         pyplot.subplot(2, 1, 1)
6         pyplot.title('Cross Entropy Loss')
7         pyplot.plot(histories[i].history['loss'], color='blue', label='train')
8         pyplot.plot(histories[i].history['val_loss'], color='orange', label='test')
9         # plot accuracy
10        pyplot.subplot(2, 1, 2)
11        pyplot.title('Classification Accuracy')
12        pyplot.plot(histories[i].history['accuracy'], color='blue', label='train')
13        pyplot.plot(histories[i].history['val_accuracy'], color='orange', label='test')
14    pyplot.show()

```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

Email Address

[START MY EMAIL COURSE](#)

Next, the classification accuracy scores collected during each fold can be summarized by calculating the mean and standard deviation. This provides an estimate of the average expected performance of the model trained on this dataset, with an estimate of the average variance in the mean. We will also summarize the distribution of scores by creating and showing a box and whisker plot.

The `summarize_performance()` function below implements this for a given list of scores collected during model evaluation.

```

1 # summarize model performance
2 def summarize_performance(scores):
3     # print summary
4     print('Accuracy: mean=%.3f std=%.3f, n=%d' % (mean(scores)*100, std(scores)*100, len(scores)))
5     # box and whisker plots of results
6     pyplot.boxplot(scores)
7     pyplot.show()

```

Complete Example

We need a function that will drive the test harness.

[Start Machine Learning](#)

This involves calling all of the define functions.

```

1 # run the test harness for evaluating a model
2 def run_test_harness():
3     # load dataset
4     trainX, trainY, testX, testY = load_dataset()
5     # prepare pixel data
6     trainX, testX = prep_pixels(trainX, testX)
7     # evaluate model
8     scores, histories = evaluate_model(trainX, trainY)
9     # learning curves
10    summarize_diagnostics(histories)
11    # summarize estimated performance
12    summarize_performance(scores)

```

We now have everything we need; the complete code to build and evaluate a CNN model on the MNIST dataset is listed below.

```

1 # baseline cnn model for mnist
2 from numpy import mean
3 from numpy import std
4 from matplotlib import pyplot
5 from sklearn.model_selection import KFold
6 from keras.datasets import mnist
7 from keras.utils import to_categorical
8 from keras.models import Sequential
9 from keras.layers import Conv2D
10 from keras.layers import MaxPooling2D
11 from keras.layers import Dense
12 from keras.layers import Flatten
13 from keras.optimizers import SGD
14
15 # load train and test dataset
16 def load_dataset():
17     # load dataset
18     (trainX, trainY), (testX, testY) = mnist.load_data()
19     # reshape dataset to have a single channel
20     trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
21     testX = testX.reshape((testX.shape[0], 28, 28, 1))
22     # one hot encode target values
23     trainY = to_categorical(trainY)
24     testY = to_categorical(testY)
25     return trainX, trainY, testX, testY
26
27 # scale pixels
28 def prep_pixels(train, test):
29     # convert from integers to floats
30     train_norm = train.astype('float32')
31     test_norm = test.astype('float32')
32     # normalize to range 0-1
33     train_norm = train_norm / 255.0
34     test_norm = test_norm / 255.0
35     # return normalized images
36     return train_norm, test_norm
37
38 # define cnn model
39 def define_model():
40     model = Sequential()
41     model.add(Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)))
42     model.add(MaxPooling2D((2, 2)))

```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

```

43     model.add(Flatten())
44     model.add(Dense(100, activation='relu', kernel_initializer='he_uniform'))
45     model.add(Dense(10, activation='softmax'))
46     # compile model
47     opt = SGD(lr=0.01, momentum=0.9)
48     model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
49     return model
50
51 # evaluate a model using k-fold cross-validation
52 def evaluate_model(dataX, dataY, n_folds=5):
53     scores, histories = list(), list()
54     # prepare cross validation
55     kfold = KFold(n_folds, shuffle=True, random_state=1)
56     # enumerate splits
57     for train_ix, test_ix in kfold.split(dataX):
58         # define model
59         model = define_model()
60         # select rows for train and test
61         trainX, trainY, testX, testY = dataX
62         # fit model
63         history = model.fit(trainX, trainY,
64             # evaluate model
65             -, acc = model.evaluate(testX, testY)
66             print('> %.3f' % (acc * 100.0))
67             # stores scores
68             scores.append(acc)
69             histories.append(history)
70     return scores, histories
71
72 # plot diagnostic learning curves
73 def summarize_diagnostics(histories):
74     for i in range(len(histories)):
75         # plot loss
76         pyplot.subplot(2, 1, 1)
77         pyplot.title('Cross Entropy Loss')
78         pyplot.plot(histories[i].history['loss'], color='blue', label='train')
79         pyplot.plot(histories[i].history['val_loss'], color='orange', label='test')
80         # plot accuracy
81         pyplot.subplot(2, 1, 2)
82         pyplot.title('Classification Accuracy')
83         pyplot.plot(histories[i].history['accuracy'], color='blue', label='train')
84         pyplot.plot(histories[i].history['val_accuracy'], color='orange', label='test')
85     pyplot.show()
86
87 # summarize model performance
88 def summarize_performance(scores):
89     # print summary
90     print('Accuracy: mean=%.3f std=%.3f, n=%d' % (mean(scores)*100, std(scores)*100, len(scores)))
91     # box and whisker plots of results
92     pyplot.boxplot(scores)
93     pyplot.show()
94
95 # run the test harness for evaluating a model
96 def run_test_harness():
97     # load dataset
98     trainX, trainY, testX, testY = load_dataset()
99     # prepare pixel data
100    trainX, testX = prep_pixels(trainX, testX)
101    # evaluate model
102    scores, histories = evaluate_model(trainX, trainY)
103    # learning curves
104    summarize_diagnostics(histories)
105    # summarize estimated performance
106    summarize_performance(scores)
107
108    ...

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

[Start Machine Learning](#)

```
108 # entry point, run the test harness  
109 run_test_harness()
```

Running the example prints the classification accuracy for each fold of the cross-validation process. This is helpful to get an idea that the model evaluation is progressing.

Note: Your [results may vary](#) given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

We can see two cases where the model achieves perfect skill and one case where it achieved lower than 98% accuracy. These are good results.

```
1 > 98.550  
2 > 98.600  
3 > 98.642  
4 > 98.850  
5 > 98.742
```

Next, a diagnostic plot is shown, giving insight into the training progress.

In this case, we can see that the model generally achieves convergence quickly, with the training loss converging. There is no obvious sign of over- or underfitting.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

[Start Machine Learning](#)

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loss and Accuracy Learning Curves for the Baseline Model During k-Fold Cross-Validation

Next, a summary of the model performance is calculated.

We can see in this case, the model has an estimated skill of about 98.6%, which is reasonable.

```
1 Accuracy: mean=98.677 std=0.107, n=5
```

Finally, a box and whisker plot is created to summarize the distribution of accuracy scores.

Start Machine Learning

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Box and Whisker Plot of Accuracy Scores for the Baseline Model Evaluated Using k-Fold Cross-Validation

We now have a robust test harness and a well-performing baseline model.

How to Develop an Improved Model

There are many ways that we might explore improvements to the baseline model.

We will look at areas of model configuration that often result in an improvement, so-called low-hanging fruit. The first is a change to the learning algorithm, and the second is an increase in the depth of the model.

Improvement to Learning

There are many aspects of the learning algorithm that can be explored for improvement.

Perhaps the point of biggest leverage is the learning rate, such as evaluating the impact that smaller or larger values of the learning rate may have, as well as schedules that change the learning rate during training.

Start Machine Learning

Another approach that can rapidly accelerate the learning of a model and can result in large performance improvements is batch normalization. We will evaluate the effect that batch normalization has on our baseline model.

Batch normalization can be used after convolutional and fully connected layers. It has the effect of changing the distribution of the output of the layer, specifically by standardizing the outputs. This has the effect of stabilizing and accelerating the learning process.

We can update the model definition to use batch normalization after the activation function for the convolutional and dense layers of our baseline model. The updated version of `define_model()` function with batch normalization is listed below.

```

1 # define cnn model
2 def define_model():
3     model = Sequential()
4     model.add(Conv2D(32, (3, 3), activation='relu'))
5     model.add(BatchNormalization())
6     model.add(MaxPooling2D((2, 2)))
7     model.add(Flatten())
8     model.add(Dense(100, activation='relu'))
9     model.add(BatchNormalization())
10    model.add(Dense(10, activation='softmax'))
11    # compile model
12    opt = SGD(lr=0.01, momentum=0.9)
13    model.compile(optimizer=opt, loss='categorical_crossentropy')
14    return model

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

The complete code listing with this change is provided below.

```

1 # cnn model with batch normalization for mnist
2 from numpy import mean
3 from numpy import std
4 from matplotlib import pyplot
5 from sklearn.model_selection import KFold
6 from keras.datasets import mnist
7 from keras.utils import to_categorical
8 from keras.models import Sequential
9 from keras.layers import Conv2D
10 from keras.layers import MaxPooling2D
11 from keras.layers import Dense
12 from keras.layers import Flatten
13 from keras.optimizers import SGD
14 from keras.layers import BatchNormalization
15
16 # load train and test dataset
17 def load_dataset():
18     # load dataset
19     (trainX, trainY), (testX, testY) = mnist.load_data()
20     # reshape dataset to have a single channel
21     trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
22     testX = testX.reshape((testX.shape[0], 28, 28, 1))
23     # one hot encode target values
24     trainY = to_categorical(trainY)
25     testY = to_categorical(testY)
26     return trainX, trainY, testX, testY
27

```

[Start Machine Learning](#)

```

28 # scale pixels
29 def prep_pixels(train, test):
30     # convert from integers to floats
31     train_norm = train.astype('float32')
32     test_norm = test.astype('float32')
33     # normalize to range 0-1
34     train_norm = train_norm / 255.0
35     test_norm = test_norm / 255.0
36     # return normalized images
37     return train_norm, test_norm
38
39 # define cnn model
40 def define_model():
41     model = Sequential()
42     model.add(Conv2D(32, (3, 3), activation='relu', kernel_initializer='he_uniform', input_shape=(28, 28, 1)))
43     model.add(BatchNormalization())
44     model.add(MaxPooling2D((2, 2)))
45     model.add(Flatten())
46     model.add(Dense(100, activation='relu'))
47     model.add(BatchNormalization())
48     model.add(Dense(10, activation='softmax'))
49     # compile model
50     opt = SGD(lr=0.01, momentum=0.9)
51     model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
52     return model
53
54 # evaluate a model using k-fold cross-validation
55 def evaluate_model(dataX, dataY, n_folds=5):
56     scores, histories = list(), list()
57     # prepare cross validation
58     kfold = KFold(n_folds, shuffle=True, random_state=seed)
59     # enumerate splits
60     for train_ix, test_ix in kfold.split(dataX):
61         # define model
62         model = define_model()
63         # select rows for train and test
64         trainX, trainY, testX, testY = dataX[train_ix], dataY[train_ix], dataX[test_ix], dataY[test_ix]
65         # fit model
66         history = model.fit(trainX, trainY, epochs=10, batch_size=32, validation_data=(testX, testY))
67         # evaluate model
68         _, acc = model.evaluate(testX, testY, verbose=0)
69         print('> %.3f' % (acc * 100.0))
70         # stores scores
71         scores.append(acc)
72         histories.append(history)
73     return scores, histories
74
75 # plot diagnostic learning curves
76 def summarize_diagnostics(histories):
77     for i in range(len(histories)):
78         # plot loss
79         pyplot.subplot(2, 1, 1)
80         pyplot.title('Cross Entropy Loss')
81         pyplot.plot(histories[i].history['loss'], color='blue', label='train')
82         pyplot.plot(histories[i].history['val_loss'], color='orange', label='test')
83         # plot accuracy
84         pyplot.subplot(2, 1, 2)
85         pyplot.title('Classification Accuracy')
86         pyplot.plot(histories[i].history['accuracy'], color='blue', label='train')
87         pyplot.plot(histories[i].history['val_accuracy'], color='orange', label='test')
88     pyplot.show()
89
90 # summarize model performance
91 def summarize_performance(scores):
92     # print summary

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

[Start Machine Learning](#)

```

93     print('Accuracy: mean=%.3f std=%.3f, n=%d' % (mean(scores)**100, std(scores)**100, len(scores)))
94     # box and whisker plots of results
95     pyplot.boxplot(scores)
96     pyplot.show()
97
98 # run the test harness for evaluating a model
99 def run_test_harness():
100     # load dataset
101     trainX, trainY, testX, testY = load_dataset()
102     # prepare pixel data
103     trainX, testX = prep_pixels(trainX, testX)
104     # evaluate model
105     scores, histories = evaluate_model(trainX, trainY)
106     # learning curves
107     summarize_diagnostics(histories)
108     # summarize estimated performance
109     summarize_performance(scores)
110
111 # entry point, run the test harness
112 run_test_harness()

```

Running the example again reports model performance.

Note: Your [results may vary](#) given the stochastic nature of the training process and differences in numerical precision. Consider running the example multiple times to get a feel for the outcome.

We can see perhaps a small drop in model performance from the first validation fold to the last validation fold.

```

1 > 98.475
2 > 98.608
3 > 98.683
4 > 98.783
5 > 98.667

```

A plot of the learning curves is created, in this case showing that the speed of learning (improvement over epochs) does not appear to be different from the baseline model.

The plots suggest that batch normalization, at least as implemented in this case, does not offer any benefit.

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

Start Machine Learning ×

You can master applied Machine Learning
without math or fancy degrees.
Find out how in this *free* and *practical* course.

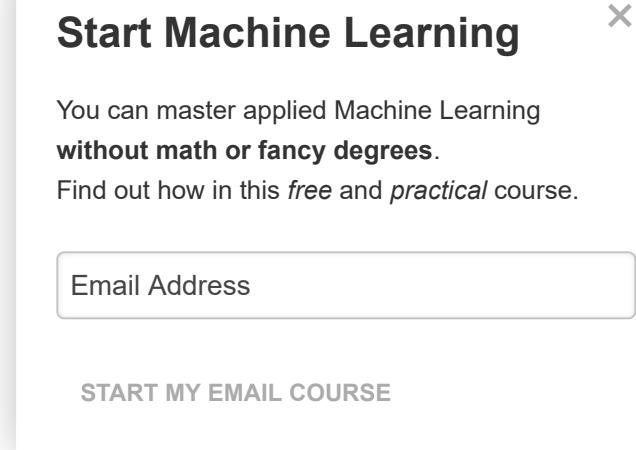
START MY EMAIL COURSE

Loss and Accuracy Learning Curves for the BatchNormalization Model During k-Fold Cross-Validation

Next, the estimated performance of the model is presented, showing performance with a slight decrease in the mean accuracy of the model: 98.643 as compared to 98.677 with the baseline model.

```
1 Accuracy: mean=98.643 std=0.101, n=5
```

Start Machine Learning



Box and Whisker Plot of Accuracy Scores for the BatchNormalization Model Evaluated Using k-Fold Cross-Validation

Increase in Model Depth

There are many ways to change the model configuration in order to explore improvements over the baseline model.

Two common approaches involve changing the [capacity](#) of the feature extraction part of the model or changing the capacity or function of the classifier part of the model. Perhaps the point of biggest influence is a change to the feature extractor.

We can increase the depth of the feature extractor part of the model, following a [VGG-like pattern](#) of adding more convolutional and pooling layers with the same sized filter, while increasing the number of filters. In this case, we will add a double convolutional layer with 64 filters each, followed by another max pooling layer.

The updated version of the `define_model()` function with this change is listed below.

```
1 # define cnn model
2 def define_model():
3     ...
```

Start Machine Learning

```

3     model = Sequential()
4     model.add(Conv2D(32, (3, 3), activation='relu', kernel_initializer='he_uniform', input_shap
5     model.add(MaxPooling2D((2, 2)))
6     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
7     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
8     model.add(MaxPooling2D((2, 2)))
9     model.add(Flatten())
10    model.add(Dense(100, activation='relu', kernel_initializer='he_uniform'))
11    model.add(Dense(10, activation='softmax'))
12    # compile model
13    opt = SGD(lr=0.01, momentum=0.9)
14    model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
15    return model

```

For completeness, the entire code listing, including this change, is provided below.

```

1 # deeper cnn model for mnist
2 from numpy import mean
3 from numpy import std
4 from matplotlib import pyplot
5 from sklearn.model_selection import KFold
6 from keras.datasets import mnist
7 from keras.utils import to_categorical
8 from keras.models import Sequential
9 from keras.layers import Conv2D
10 from keras.layers import MaxPooling2D
11 from keras.layers import Dense
12 from keras.layers import Flatten
13 from keras.optimizers import SGD
14
15 # load train and test dataset
16 def load_dataset():
17     # load dataset
18     (trainX, trainY), (testX, testY) = mnist.load_data()
19     # reshape dataset to have a single channel
20     trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
21     testX = testX.reshape((testX.shape[0], 28, 28, 1))
22     # one hot encode target values
23     trainY = to_categorical(trainY)
24     testY = to_categorical(testY)
25     return trainX, trainY, testX, testY
26
27 # scale pixels
28 def prep_pixels(train, test):
29     # convert from integers to floats
30     train_norm = train.astype('float32')
31     test_norm = test.astype('float32')
32     # normalize to range 0-1
33     train_norm = train_norm / 255.0
34     test_norm = test_norm / 255.0
35     # return normalized images
36     return train_norm, test_norm
37
38 # define cnn model
39 def define_model():
40     model = Sequential()
41     model.add(Conv2D(32, (3, 3), activation='relu', kernel_initializer='he_uniform', input_shap
42     model.add(MaxPooling2D((2, 2)))
43     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
44     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
45     model.add(MaxPooling2D((2, 2)))
46     model.add(Flatten())
47     model.add(Dense(100, activation='relu',

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

[Start Machine Learning](#)

```

48     model.add(Dense(10, activation='softmax'))
49     # compile model
50     opt = SGD(lr=0.01, momentum=0.9)
51     model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
52     return model
53
54 # evaluate a model using k-fold cross-validation
55 def evaluate_model(dataX, dataY, n_folds=5):
56     scores, histories = list(), list()
57     # prepare cross validation
58     kfold = KFold(n_folds, shuffle=True, random_state=1)
59     # enumerate splits
60     for train_ix, test_ix in kfold.split(dataX):
61         # define model
62         model = define_model()
63         # select rows for train and test
64         trainX, trainY, testX, testY = dataX[train_ix], dataY[train_ix], dataX[test_ix], dataY[test_ix]
65         # fit model
66         history = model.fit(trainX, trainY,
67             # evaluate model
68             _, acc = model.evaluate(testX, testY)
69             print('> %.3f' % (acc * 100.0))
70         # stores scores
71         scores.append(acc)
72         histories.append(history)
73     return scores, histories
74
75 # plot diagnostic learning curves
76 def summarize_diagnostics(histories):
77     for i in range(len(histories)):
78         # plot loss
79         pyplot.subplot(2, 1, 1)
80         pyplot.title('Cross Entropy Loss')
81         pyplot.plot(histories[i].history['loss'], color='blue', label='train')
82         pyplot.plot(histories[i].history['val_loss'], color='orange', label='test')
83         # plot accuracy
84         pyplot.subplot(2, 1, 2)
85         pyplot.title('Classification Accuracy')
86         pyplot.plot(histories[i].history['accuracy'], color='blue', label='train')
87         pyplot.plot(histories[i].history['val_accuracy'], color='orange', label='test')
88     pyplot.show()
89
90 # summarize model performance
91 def summarize_performance(scores):
92     # print summary
93     print('Accuracy: mean=%.3f std=%.3f, n=%d' % (mean(scores)*100, std(scores)*100, len(scores)))
94     # box and whisker plots of results
95     pyplot.boxplot(scores)
96     pyplot.show()
97
98 # run the test harness for evaluating a model
99 def run_test_harness():
100     # load dataset
101     trainX, trainY, testX, testY = load_dataset()
102     # prepare pixel data
103     trainX, testX = prep_pixels(trainX, testX)
104     # evaluate model
105     scores, histories = evaluate_model(trainX, trainY)
106     # learning curves
107     summarize_diagnostics(histories)
108     # summarize estimated performance
109     summarize_performance(scores)
110
111 # entry point, run the test harness
112 run_test_harness()

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

[Start Machine Learning](#)

Running the example reports model performance for each fold of the cross-validation process.

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

The per-fold scores may suggest some improvement over the baseline.

```
1 > 99.058  
2 > 99.042  
3 > 98.883  
4 > 99.192  
5 > 99.133
```

A plot of the learning curves is created, in this case showing a good fit to the training problem, with no clear signs of overfitting. The plots may be helpful.

Start Machine Learning X

You can master applied Machine Learning without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Loss and Accuracy Learning Curves for the Deeper Model During k-Fold Cross-Validation

Start Machine Learning

Next, the estimated performance of the model is presented, showing a small improvement in performance as compared to the baseline from 98.677 to 99.062, with a small drop in the standard deviation as well.

1 Accuracy: mean=99.062 std=0.104, n=5

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Box and Whisker Plot of Accuracy Scores for the Deeper Model Evaluated Using k-Fold Cross-Validation

How to Finalize the Model and Make Predictions

The process of model improvement may continue for as long as we have ideas and the time and resources to test them out.

At some point, a final model configuration must be chosen and adopted. In this case, we will choose the deeper model as our final model.

First, we will finalize our model, but fitting a model on the entire training dataset and saving the model to file for later use. We will then load the model and evaluate its performance on the hold out test dataset to get an idea of how well the chosen model actually performs in practice. Finally, we will use the saved model to make a prediction on a single image.

[Start Machine Learning](#)

Save Final Model

A final model is typically fit on all available data, such as the combination of all train and test dataset.

In this tutorial, we are intentionally holding back a test dataset so that we can estimate the performance of the final model, which can be a good idea in practice. As such, we will fit our model on the training dataset only.

```
1 # fit model
2 model.fit(trainX, trainY, epochs=10, batch_size=32, verbose=0)
```

Once fit, we can save the final model to an H5 file by calling the `save()` function on the model and pass in the chosen filename.

```
1 # save model
2 model.save('final_model.h5')
```

Note, saving and loading a Keras model requires that

The complete example of fitting the final deep model can be found below.

```
1 # save the final model to file
2 from keras.datasets import mnist
3 from keras.utils import to_categorical
4 from keras.models import Sequential
5 from keras.layers import Conv2D
6 from keras.layers import MaxPooling2D
7 from keras.layers import Dense
8 from keras.layers import Flatten
9 from keras.optimizers import SGD
10
11 # load train and test dataset
12 def load_dataset():
13     # load dataset
14     (trainX, trainY), (testX, testY) = mnist.load_data()
15     # reshape dataset to have a single channel
16     trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
17     testX = testX.reshape((testX.shape[0], 28, 28, 1))
18     # one hot encode target values
19     trainY = to_categorical(trainY)
20     testY = to_categorical(testY)
21     return trainX, trainY, testX, testY
22
23 # scale pixels
24 def prep_pixels(train, test):
25     # convert from integers to floats
26     train_norm = train.astype('float32')
27     test_norm = test.astype('float32')
28     # normalize to range 0-1
29     train_norm = train_norm / 255.0
30     test_norm = test_norm / 255.0
31     # return normalized images
32     return train_norm, test_norm
33
```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

[Start Machine Learning](#)

```

54 # define cnn model
55 def define_model():
56     model = Sequential()
57     model.add(Conv2D(32, (3, 3), activation='relu', kernel_initializer='he_uniform', input_shape=(28, 28, 1)))
58     model.add(MaxPooling2D((2, 2)))
59     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
60     model.add(Conv2D(64, (3, 3), activation='relu', kernel_initializer='he_uniform'))
61     model.add(MaxPooling2D((2, 2)))
62     model.add(Flatten())
63     model.add(Dense(100, activation='relu', kernel_initializer='he_uniform'))
64     model.add(Dense(10, activation='softmax'))
65     # compile model
66     opt = SGD(lr=0.01, momentum=0.9)
67     model.compile(optimizer=opt, loss='categorical_crossentropy', metrics=['accuracy'])
68     return model
69
70 # run the test harness for evaluating a model
71 def run_test_harness():
72     # load dataset
73     trainX, trainY, testX, testY = load_data()
74     # prepare pixel data
75     trainX, testX = prep_pixels(trainX, testX)
76     # define model
77     model = define_model()
78     # fit model
79     model.fit(trainX, trainY, epochs=10, batch_size=64)
80     # save model
81     model.save('final_model.h5')
82
83 # entry point, run the test harness
84 run_test_harness()

```

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

After running this example, you will now have a 1.2-megabyte file named "final_model.h5" in your current working directory.

Evaluate Final Model

We can now load the final model and evaluate it on the hold out test dataset.

This is something we might do if we were interested in presenting the performance of the chosen model to project stakeholders.

The model can be loaded via the `load_model()` function.

The complete example of loading the saved model and evaluating it on the test dataset is listed below.

```

1 # evaluate the deep model on the test dataset
2 from keras.datasets import mnist
3 from keras.models import load_model
4 from keras.utils import to_categorical
5
6 # load train and test dataset
7 def load_dataset():
8     # load dataset
9     (trainX, trainY), (testX, testY) = mnist.load_data()
10    # reshape dataset to have a single channel
11    trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
12    testX = testX.reshape((testX.shape[0], 28, 28, 1))
13    # one hot encode target values

```

[Start Machine Learning](#)

```

14     trainY = to_categorical(trainY)
15     testY = to_categorical(testY)
16     return trainX, trainY, testX, testY
17
18 # scale pixels
19 def prep_pixels(train, test):
20     # convert from integers to floats
21     train_norm = train.astype('float32')
22     test_norm = test.astype('float32')
23     # normalize to range 0-1
24     train_norm = train_norm / 255.0
25     test_norm = test_norm / 255.0
26     # return normalized images
27     return train_norm, test_norm
28
29 # run the test harness for evaluating a model
30 def run_test_harness():
31     # load dataset
32     trainX, trainY, testX, testY = load_datal...
33     # prepare pixel data
34     trainX, testX = prep_pixels(trainX, testX)
35     # load model
36     model = load_model('final_model.h5')
37     # evaluate model on test dataset
38     _, acc = model.evaluate(testX, testY, ver...
39     print('> %.3f' % (acc * 100.0))
40
41 # entry point, run the test harness
42 run_test_harness()

```

Running the example loads the saved model and eval

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

Note: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

The classification accuracy for the model on the test dataset is calculated and printed. In this case, we can see that the model achieved an accuracy of 99.090%, or just less than 1%, which is not bad at all and reasonably close to the estimated 99.753% with a standard deviation of about half a percent (e.g. 99% of scores).

1 > 99.090

Make Prediction

We can use our saved model to make a prediction on new images.

The model assumes that new images are grayscale, that they have been aligned so that one image contains one centered handwritten digit, and that the size of the image is square with the size 28×28 pixels.

Below is an image extracted from the MNIST test dataset. You can save it in your current working directory with the filename ‘sample_image.png’.

[Start Machine Learning](#)

- Download the sample image ([sample_image.png](#))

We will pretend this is an entirely new and unseen image. We might use our saved model to predict the integer that it represents.

First, we can load the image, force it to be in grayscale, and the loaded image can then be resized to have a single channel. The `load_image()` function implements this and will return the loaded image ready for classification.

Importantly, the pixel values are prepared in the same way as the pixel values were prepared for the training dataset when fitting the final model, in this case, normalized.

```

1 # load and prepare the image
2 def load_image(filename):
3     # load the image
4     img = load_img(filename, grayscale=True, target_size=(28, 28))
5     # convert to array
6     img = img_to_array(img)
7     # reshape into a single sample with 1 channel
8     img = img.reshape(1, 28, 28, 1)
9     # prepare pixel data
10    img = img.astype('float32')
11    img = img / 255.0
12    return img

```

Next, we can load the model as in the previous section and call the `predict_classes()` function to predict the digit that the image represents.

```

1 # predict the class
2 digit = model.predict_classes(img)

```

The complete example is listed below.

Start Machine Learning

You can master applied Machine Learning without **math or fancy degrees**.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

```

1 # make a prediction for a new image.
2 from keras.preprocessing.image import load_img
3 from keras.preprocessing.image import img_to_array
4 from keras.models import load_model
5
6 # load and prepare the image
7 def load_image(filename):
8     # load the image
9     img = load_img(filename, grayscale=True, target_size=(28, 28))
10    # convert to array
11    img = img_to_array(img)
12    # reshape into a single sample with 1 channel
13    img = img.reshape(1, 28, 28, 1)
14    # prepare pixel data
15    img = img.astype('float32')
16    img = img / 255.0
17    return img
18
19 # load an image and predict the class
20 def run_example():
21     # load the image
22     img = load_image('sample_image.png')
23     # load model
24     model = load_model('final_model.h5')
25     # predict the class
26     digit = model.predict_classes(img)
27     print(digit[0])
28
29 # entry point, run the example
30 run_example()

```

Running the example first loads and prepares the image. The loaded image represents the digit ‘7’.

1 7

Extensions

This section lists some ideas for extending the tutorial that you may wish to explore.

- **Tune Pixel Scaling.** Explore how alternate pixel scaling methods impact model performance as compared to the baseline model, including centering and standardization.
- **Tune the Learning Rate.** Explore how different learning rates impact the model performance as compared to the baseline model, such as 0.001 and 0.0001.
- **Tune Model Depth.** Explore how adding more layers to the model impact the model performance as compared to the baseline model, such as another block of convolutional and pooling layers or another dense layer in the classifier part of the model.

If you explore any of these extensions, I’d love to know.

Post your findings in the comments below.

Further Reading

This section provides more resources on the topic if you’re interested.

[Start Machine Learning](#)

APIs

- Keras Datasets API
- Keras Datasets Code
- sklearn.model_selection.KFold API

Articles

- MNIST database, Wikipedia.
- Classification datasets results, What is the class of this image?

Summary

In this tutorial, you discovered how to develop a convolutional neural network for handwritten digit classification from scratch.

Specifically, you learned:

- How to develop a test harness to develop a robust baseline model with 90% performance for a classification task.
- How to explore extensions to a baseline model to increase performance.
- How to develop a finalized model, evaluate the performance, and make predictions on new images.

Do you have any questions?

Ask your questions in the comments below and I will do my best to answer.

Start Machine Learning

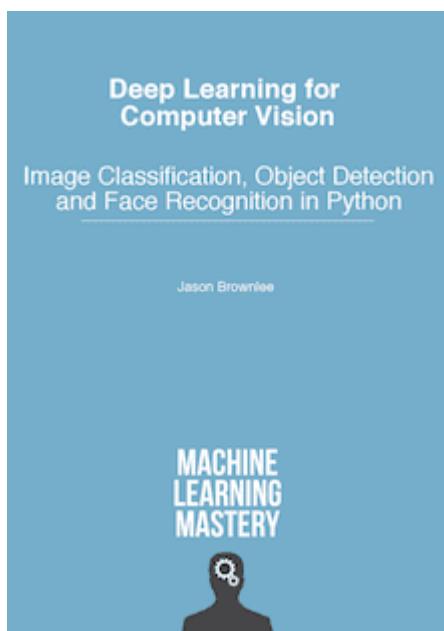
You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Develop Deep Learning Models for Vision Today!



Develop Your Own Vision Models in Minutes

...with just a few lines of python code

Discover how in my new Ebook:
Deep Learning for Computer Vision

It provides **self-study tutorials** on topics like:
classification, object detection (yolo and rcnn), face recognition (vggface and facenet), data preparation and much more...

Finally Bring Deep Learning to your Vision Projects

Skip the Academics. Just Results.

SEE WHAT'S INSIDE

Start Machine Learning

[Tweet](#)[Share](#)[Share](#)

About Jason Brownlee

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

[View all posts by Jason Brownlee →](#)

◀ How to Visualize Filters and Feature Maps in Convolutional Neural Networks

140 Responses to *How to Develop a CNN for MNIST Handwritten Digit Classification*

SHAHEEN ALHIRMIZY May 10, 2019 at 4:32 pm #

Dear Jason Thank You Very Much For your w
for Edge detection from scratch using real images.

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

Jason Brownlee May 11, 2019 at 6:05 am #

[REPLY ↗](#)

Great suggestion, thanks.

SHAHEEN ALHIRMIZY May 10, 2019 at 4:47 pm #

[REPLY ↗](#)

I have two questions :

first suppose I have images sizes equal to 7611 x 7811 how to deal with this big size images in CNN models.

second question about if our images different sizes not same size.

Jason Brownlee May 11, 2019 at 6:06 am #

[REPLY ↗](#)

I recommend reducing the size of your images first, before modeling, e.g. less than 1000 pixels,
even less than 500 pixels if you can.

[Start Machine Learning](#)

I then recommend normalizing images to the same size.

vinay September 30, 2019 at 11:47 am #

REPLY ↗

to reduce the size of the images, which approach is better. Using machine learning or normal image compression methods.

Jason Brownlee September 30, 2019 at 2:26 pm #

REPLY ↗

Use image compression algorithms.

SHAHEEN ALHIRMIZY May 11, 2019 at 4:13 pm #

when we reducing the size of images that me:

Jason Brownlee May 12, 2019 at 6:38 am #

Often it does.

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Taranpreetkaur October 12, 2019 at 10:57 pm #

REPLY ↗

How u make this project please help us ... We can't understand how we make this project ...
Taranpreetkaur

Jason Brownlee October 13, 2019 at 8:30 am #

REPLY ↗

What problem are you having exactly?

Ajoy October 20, 2019 at 4:08 am #

REPLY ↗

Original error was: No module named _multiarray_umath, keras .imageprocessing
Importing the numpy c-extensions failed.

Any ideas please

Start Machine Learning

Jason Brownlee October 20, 2019 at 6:23 am #

REPLY ↗

Sorry to hear that.

Perhaps try checking that your version of Keras and TensorFlow are up to date?

Ajoy October 20, 2019 at 7:18 pm #

REPLY ↗

Sorry dont worry, I somehow got to resolve it. I think I reinstalled stuff.

I have another problem though. How do I add data for 17? I am not sure what to do. pelase. I think I need the advise for both how to add as well as how to split. Thanks

Jason Brownlee October 21, 2019 at 6:17 am #

Happy to hear that.

Perhaps you can refit your model as you get access to the new data.

Start Machine Learning

You can master applied Machine Learning without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Ajoy October 20, 2019 at 7:51 pm #

REPLY ↗

Jason you know what, for “17” it returned “7”. Thats actually good, but if you could please tip me as to how can I turn this into “Awesome.” Can I divide the pixel data into 2 dividions based on whitespace or what would be the best way please to recognize “17”. Thanks

Jason Brownlee October 21, 2019 at 6:17 am #

REPLY ↗

I believe the model supports one character at a time, try splitting multiple characters up into single character input.

Jubayed November 26, 2019 at 10:13 pm #

REPLY ↗

Great!! I enjoyed very well.

Jason Brownlee November 27, 2019 at 6:05 pm #

Start Machine Learning

Thanks!

Victor Johns November 27, 2019 at 3:43 am #

REPLY ↗

Jason,

Thanks much for this tutorial. I had a question on the structure of the CNN. In this (and other descriptions of this problem for the MNIST digit data set) a common structure seems to be that the number of CNN filters is 32. Where each filter is a 5×5 of stride 1. However (and I apologize if this is explained somewhere) I cannot seem to find why 32 filters are chosen. i.e. Why not 16 or 64 (or 14 or 28). Is it possible to explain why this is so.

Regards

Victor

Jason Brownlee November 27, 2019 at 6:13 pm #

It is arbitrary.

Try different numbers of filters and compare the results.

Start Machine Learning



You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Samuel December 11, 2019 at 8:36 am #

REPLY ↗

Jason,

Thank you for the detailed explanation of this. I wanted to run your code and follow along to see how it worked. When I tried to run the complete example for the baseline cnn model for mnist, I get the following error:

```
KeyError Traceback (most recent call last)
in
107
108 # entry point, run the test harness
-> 109 run_test_harness()

in run_test_harness()
102 scores, histories = evaluate_model(model, trainX, trainY)
103 # learning curves
-> 104 summarize_diagnostics(histories)
105 # summarize estimated performance
106 summarize_performance(scores)

in summarize_diagnostics(histories)
79 pyplot.subplot(212)
80 pyplot.title('Classification Accuracy')
-> 81 pyplot.plot(histories[i].history['acc'], color='blue')
```

Start Machine Learning

```
82 pyplot.plot(histories[i].history['val_acc'], color='orange', label='test')
83 pyplot.show()
```

KeyError: 'acc'

Any ideas what it means?

Jason Brownlee December 11, 2019 at 1:41 pm #

REPLY ↗

Yes, the API has changed.

I have updated the example.

Thanks for letting me know!

Xenon December 12, 2019 at 5:16 am #

please how can i Train a simple convolutional
that's what i want to see.

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

Jason Brownlee December 12, 2019 at 6:33 am #

See the above tutorial for an example.

Vivek December 23, 2019 at 10:30 pm #

REPLY ↗

I want to create a image exactly same as mnist using mspaint for example and test it using the
mnist evaluation model using predict_classes API.

Is it possible without using python API.

Jason Brownlee December 24, 2019 at 6:42 am #

REPLY ↗

I don't see why not.

Alexander Soare December 27, 2019 at 8:00 am #

REPLY ↗

Hi Jason,

Thanks for the excellent tutorial! I have two questions about

Start Machine Learning

introduction" article you say that the model is discarded each time we switch the hold out set. So here are my questions

- 1) Here it looks like you don't discard the model between iterations of the k-fold cross validation. Why?
- 2) The model's accuracy gets better each time you do an iteration. Why is taking a simple average fair then? To take it to the extreme: What if you only trained on 1 epoch? Then you'd get 94% on your first fold, then 98.5% on your second fold. And so on, yielding a much poorer average and a huge standard deviation. On the back of this question, I'd also ask why the standard deviation is a fair representation of the model under these circumstances.

Thanks again for the great material!

Alex

Jason Brownlee December 28, 2019 at 7:38

You're welcome.

That looks like an error. I will update the tutorial.

Thanks for pointing it out!

Start Machine Learning

You can master applied Machine Learning without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Alexander Soare December 29, 2019 at 5:43 am #

Thanks for picking up my previous question Jason. I have another. You say at some point that there are no clear signs of overfitting. But it looks like the orange lines for the validation accuracy are staying relatively low compared to the blue lines for the testing accuracy. Is this not an indicator of overfitting? Or am I misinterpreting the charts?

Jason Brownlee December 29, 2019 at 6:09 am #

REPLY ↗

Excellent question Alexander.

They are close enough that I would not classify it as overfitting, instead it is a good fit:

<https://machinelearningmastery.com/learning-curves-for-diagnosing-machine-learning-model-performance/>

Anja January 3, 2020 at 3:02 pm #

REPLY ↗

Hello Jason,

When I execute the final code I get following error: "AttributeError: module 'tensorflow' has no attribute 'get_default_graph'". I have installed Tensorflow 2.0, wi

Start Machine Learning

problem?

Thank you for your time and this amazing post!

Jason Brownlee January 4, 2020 at 8:24 am #

REPLY ↗

Confirm you have Keras 2.3 and Tensorflow 2.0 installed.

Anja January 6, 2020 at 12:51 pm #

REPLY ↗

Hello Jason,

I have Tensorflow 2.0 and Keras 2.2.4. I tried to install them on Windows 10. I tried also switching to Linux on my laptop. If you have any advice for me, I would be really grateful.

Jason Brownlee January 6, 2020 at 1:01 pm #

You can try:

```
1 pip install keras
```

Or, if you are on anaconda:

```
1 conda install -c conda-forge keras
```

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

REPLY ↗

Dmytro February 4, 2020 at 7:44 am #

REPLY ↗

Hello Jason,

Thank you for the tutorial and all the answers!

I have my own set of images, each class in a separate folder. What should I do to upload the dataset instead of writing

```
# load dataset
(trainX, trainY), (testX, testY) = mnist.load_data()
```

It would be great if you can give some advice!

Jason Brownlee February 4, 2020 at 8:02 am #

Start Machine Learning

Here is an example:

<https://machinelearningmastery.com/how-to-load-large-datasets-from-directories-for-deep-learning-with-keras/>

Dmytro February 5, 2020 at 4:00 pm #

REPLY ↗

Thank you!

Jason Brownlee February 6, 2020 at 8:18 am #

REPLY ↗

You're welcome.

SARVANI CH March 14, 2020 at 9:11 pm #

Hi Jason,

Thanks for the article!!!

Do you have any plans to implement "Information bottleneck" information among various layers in an architecture?
Does it seem simpler or a harder one for you??

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee March 15, 2020 at 6:13 am #

REPLY ↗

I have not heard of it sorry, do you have a link?

SWETHA March 20, 2020 at 1:05 pm #

REPLY ↗

```
def load_image(filename):
    # load the image
    img = load_img(filename, grayscale=True, target_size=(28, 28))
    # convert to array
    img = img_to_array(img)
    # reshape into a single sample with 1 channel
    img = img.reshape(1, 28, 28, 1)
    # prepare pixel data
    img = img.astype('float32')
    img = img / 255.0
    return img
```

Start Machine Learning

what does the 'filename' mean in the above code?

Jason Brownlee March 20, 2020 at 1:20 pm #

REPLY ↗

It's an argument to the function – the name of the file you want to load.

SWETHA March 20, 2020 at 2:17 pm #

REPLY ↗

Did you mean the saved model's filename?

Jason Brownlee March 21, 2020 at 1:00 pm #

No, the function loads an image as

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

SWETHA March 20, 2020 at 1:48 pm #

Thank you!

Jason Brownlee March 21, 2020 at 8:16 am #

REPLY ↗

You're welcome.

Elbhednam March 26, 2020 at 6:52 am #

REPLY ↗

What is the purpose of reshaping the data to a single color channel? Is it necessary for the 2D convolution step in the model definition? I ask because other digit classification examples I've looked at go straight to flattening the train/test data but those didn't have a 2D convolution layer.

Jason Brownlee March 26, 2020 at 8:07 am #

REPLY ↗

The CNN layer expects data to have a 3d shape including a channels dimension:
<https://machinelearningmastery.com/convolutional-layers-for-deep-learning-neural-networks/>

Start Machine Learning

Edward April 13, 2020 at 5:00 am #

REPLY ↗

If one wanted to create a new MNIST how would they set the data up? In these examples we get to use curated datasets from a dispensary but that does not explain how to set up your data to run the model. For example what if I wanted to create my own MNIST to determine if the picture is a poker card, how would I go about setting the pictures and telling the machine what that picture is?

Jason Brownlee April 13, 2020 at 6:22 am #

REPLY ↗

Good question, see this:

<https://machinelearningmastery.com/how-to-load-labeled-data-keras/>

Amit Baghel April 15, 2020 at 3:53 pm #

hi Jason Brownlee,
i want to make model handwriting alphabet so please su

Jason Brownlee April 16, 2020 at 5:57 am #

Perhaps this will help:

<https://machinelearningmastery.com/faq/single-faq/where-can-i-get-a-dataset-on-handwriting-alphabets/>

Nawal April 17, 2020 at 6:03 pm #

REPLY ↗

Jason, Thank you so much for sharing your knowledge. I just started learning python, tensorflow and machine learning. I understood all the maths and processes but was having difficulty in coding. Your example worked first time. Brilliant.

Do you have a simpler example of CNN with just 4 hidden layers as follows:

1. Input layer of 784 nodes
2. First convolution later : 5x5x32
3. First max pooling layer
4. Second convolution layer 5x5x64
5. Second max pooling layer
6. Out put layer of 10 nodes

Jason Brownlee April 18, 2020 at 5:42 am #

Start Machine Learning

Thanks.

You can adapt the above example to this directly if you like.

Joyce May 4, 2020 at 6:03 am #

REPLY ↩

Hi Jason. Can I ask why did you reshape the data? The image from the MNIST are already grey scale which are 1 channel. I am not quite understand why you need to do this step. Also what does the trainX.shape[0]/testX.shape[0] do here?

```
# reshape dataset to have a single channel
trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
testX = testX.reshape((testX.shape[0], 28, 28, 1))
```

Jason Brownlee May 4, 2020 at 6:28 am #

Inputs to CNNs must have the channel dimension for grayscale images.

Amy May 4, 2020 at 6:30 am #

Can I ask why the hot encode is 10 binary vectors? Why choose 10?

Jason Brownlee May 4, 2020 at 6:33 am #

REPLY ↩

Because there are 10 classes.

Ruslaniv May 7, 2020 at 9:20 pm #

REPLY ↩

Hi Jason! Thank you so much, this is a great tutorial!
Just one question, i'm somewhat confused.
Are training our model twice? Once with k-fold validation and then right before saving it, with model.fit?

====

```
model = define_model()
model.fit(trainX, trainY, epochs=10, batch_size=32, verbose=0)
model.save('final_model.h5')
=====
```

Start Machine Learning

Jason Brownlee May 8, 2020 at 6:31 am #

REPLY ↗

Cross-validation will train k models and is used to estimate the performance when making predictions on new data. These models are discarded.

Once we choose a final model config, we fit a model on all data and use it to start making predictions.

Ruslaniv May 11, 2020 at 8:20 pm #

REPLY ↗

Great, that's what I thought! Thank you again for this amazing tutorial.

Jason Brownlee May 12, 2020 at 6:31 am #

You're welcome.

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Thenerd February 5, 2021 at 3:05 am #

In the k fold cross validation a new model is trained each iteration? But then what is the use of training new models in each iteration?

By model config you mean an optimal define_model? And once an optimal model is fixed we train on the entire training set? I am not able to relate how will K iterations in the K-fold reflect our understanding of what model configuration is the best.

Jason Brownlee February 5, 2021 at 5:46 am #

REPLY ↗

Yes, a new model is trained each iteration, evaluated, and discarded.

The purpose is to estimate the performance of the model configuration when making predictions on unseen data for your prediction problem.

Once we have this estimate, we can choose to use the configuration (compared to other configurations). A final model can be fit on all data and we can start making predictions for real on new data where we don't know the label.

Thenerd February 10, 2021 at 7:21 am #

Thankyou very much. I had another question. I was trying to run the final model on the unseen test data. however for each

Start Machine Learning

run this over 10 runs and average the test results, something like monte carlo simulation. but again in that case where would model = define_model() be placed? would that be placed inside the for loop for average or just be defined once before running these 10 iterations, considering this average results is on the test data. since in kfold you discard it.

Jason Brownlee February 10, 2021 at 8:13 am #

You're welcome.

Yes, you can reduce variance in prediction by fitting multiple final models and averaging their predictions:

<https://machinelearningmastery.com/hyperparameter-optimization-with-scikit-learn/>

This has a code example:

<https://machinelearningmastery.com/multiclass-classification-with-deep-learning-tutorial/>

Gohel vivek May 11, 2020 at 6:02 pm #

If i am writing this all by making a class name

AttributeError: 'mnist_classification' object has no attribute

this .astype error is coming please sir can you help me to figure this out

Start Machine Learning X

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee May 12, 2020 at 6:40 am #

REPLY ↗

Sorry to hear that, this will help:

<https://machinelearningmastery.com/faq/single-faq/why-does-the-code-in-the-tutorial-not-work-for-me>

Ramya June 2, 2020 at 7:01 am #

REPLY ↗

Hi Jason,

Thanks for this tutorial. It's very helpful.

Jason Brownlee June 2, 2020 at 7:55 am #

REPLY ↗

You're welcome, I'm happy to hear that.

Start Machine Learning

Milind Naidu June 21, 2020 at 11:28 pm #

REPLY ↩

Hi, i am trying to build a CNN fo the same problem. My output layer has 10 neurons. When i try to fit my model after One-Hot_encoding of y_train and y_test, i get an error that the labels and the logits don't match. However, without the one-hot encoding the model fits perfectly. Can you please explain?

Jason Brownlee June 22, 2020 at 6:14 am #

REPLY ↩

Perhaps you can use the above tutorial as a starting point and adapt it for your project?

William July 19, 2020 at 8:48 pm #

This is the best article about MNIST, but I have performance in order to running these codes?"

Jason Brownlee July 20, 2020 at 6:12 am #

Thanks!

No, you can run on CPU, it might just take more minutes to run.

Sandhya July 31, 2020 at 3:58 pm #

REPLY ↩

Hi Jason

Thanks for this amazing tutorial. It will definitely help a lot of budding researchers.

I have a question

Can you please help if instead of single character i have to work on complete self taken handwritten word image (20 different classes, total around 20000 images). How to do it?

Jason Brownlee August 1, 2020 at 6:07 am #

REPLY ↩

Perhaps the letters in each word images can first be segmented, then classified.

Sandhya August 1, 2020 at 5:39 pm #

REPLY ↩

thanks for the reply but i want to do it

Start Machine Learning

Jason Brownlee August 2, 2020 at 5:39 am #

REPLY ↵

Perhaps use a CNN to read the images and an LSTM to interpret the image features and output one letter at a time.

Sandhya July 31, 2020 at 4:14 pm #

REPLY ↵

hi Jason

Thanks for this amazing tutorial.

Jason Brownlee August 1, 2020 at 6:07 am #

You're welcome.

zulfiqar ali October 18, 2020 at 9:46 am #

Thanks, Jason Brownlee for this excellent tutorial. I am your huge fan and I regularly following your all machine learning and deep earning models.I got a lot of knowledge from this site but After completion of machine learning tutorials there always raise a question in my mind that how can we apply these models in Mobile APPs(Andriod, iOS, etc) so that we get more fun and knowledge by this models. Are you suggest to me any link to have such kind of projects such as weather prediction app . Thanks again.

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee October 18, 2020 at 1:25 pm #

REPLY ↵

Thanks.

Sorry, I don't know about using models on mobile.

Zunaira Shafqat October 23, 2020 at 4:49 am #

REPLY ↵

Thanks for the tutorial Jason Brownlee.

Can I use same model but pass my own dataset (of humans) to it & use it for detecting fall or not fall by making small changes?

Thanks.

Start Machine Learning

Jason Brownlee October 23, 2020 at 6:17 am #

REPLY ↗

Sure.

Preethi November 12, 2020 at 3:28 pm #

REPLY ↗

Hi do you have any projects done based on Wireless and Mobile Network except human activity recognition using smart phone

Jason Brownlee November 13, 2020 at 6:32 pm #

I don't think so. Perhaps try the blog search

Archna November 13, 2020 at 7:01 pm #

Thanks for the detailed explanation !!
Can you pls do a tutorial on colorization without using

Thanks in advance

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

START MY EMAIL COURSE

**Jason Brownlee** November 14, 2020 at 6:30 am #

REPLY ↗

Great suggestion, thanks!

Archna November 24, 2020 at 6:31 pm #

REPLY ↗

Waiting for the same 😊

Ebdulmomen December 7, 2020 at 8:53 am #

REPLY ↗

hey jason that was a very good tutorial, i did not understand how in the predict_classes we got 7? the model had a softmax activation function as the last layer which will be a probability distribution of 10 float numbers, what i can't grasp is how did we get back 7?

Start Machine Learning

Ebdulmomen December 7, 2020 at 8:57 am #

REPLY ↗

what i mean is how this 10 probabilities turned into 1 predicted number?

Ebdulmomen December 7, 2020 at 9:09 am #

REPLY ↗

and how does it know it is 7, didn't we hot encode the Y values?

Jason Brownlee December 7, 2020 at 1:15 pm #

Via argmax:

<https://machinelearningmastery.com/argmax-in-machine-learning/>

Jason Brownlee December 7, 2020 at 1:35 pm #

Thanks.

The predict_classes() performs an argmax on the

If you are new to argmax, see this:

<https://machinelearningmastery.com/argmax-in-machine-learning/>

Start Machine Learning

X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Ebdulmomen December 7, 2020 at 6:41 pm #

REPLY ↗

thank you jason, everything is clear, really appreciate it!

Jason Brownlee December 8, 2020 at 7:40 am #

REPLY ↗

I'm happy to hear that.

Naveen December 13, 2020 at 2:43 am #

REPLY ↗

Hey is this code not compatible for Tensorflow 2? I am getting the following error.
E tensorflow/stream_executor/cuda/cuda_driver.cc:314] failed call to culinit: CUDA_ERROR_NO_DEVICE:
no CUDA-capable device is detected

Start Machine Learning

Jason Brownlee December 13, 2020 at 6:06 am #

REPLY ↗

Yes, you can learn more here:

<https://machinelearningmastery.com/faq/single-faq/do-you-support-tensorflow-2>

That error looks like a problem with your development environment.

Yash December 23, 2020 at 9:40 am #

REPLY ↗

Hello!

Thanks for the wonderful article helped me a lot with my project. I have a question though, I am trying to take an image of a handwritten digit and convert it to a number but it doesn't work with that. What should I do? I resized it to 28x28 using online tools and even tried to convert it to grayscale. Any suggestion or lead will be appreciated!

Thanks in advance!

Start Machine Learning

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee December 23, 2020 at 1:27 pm #

You will have to prepare the image in an image format with white foreground, black background, grayscale and same image size.

JG December 31, 2020 at 7:53 am #

REPLY ↗

Hi Jason,

Great tutorial as always!

I implemented your tutorial with two variants.

I load the images from digits dataset of sklearn library (it uses less images number only 1,797 and with much lower resolution 8 x 8 bits vs MNIST 70,000 and 28,28 pixels resolution. On the contrary it is much faster this dataset because it use less data features).

And also I apply for this multi-class problem other library classifier models from Sklearn such as SVC(), logisticRegression(), ExtraTreesClassifier(), XGBClassifier(), RandomForestClassifier(), etc.

My main results are:

– I got 98.4% accuracy and 1.1 of sigma for the simple CNN baseline model (not other deep layers, not other BatchNormalization()) and I got 98.8% from SVC...but only suing the 8x8 pixels resolution which is a great new!.

But when I reduce the final image to be predicted to 8 x 8 pixels, in order to apply my trained model, where I got such a great score I poorly predict the 7 digit as 9.

Start Machine Learning

My guess is when I load _image of 7 and I clip it to such smaller size of 8x8 pixels, I lost important image features in the process of cutting back the image....where the the MNIST 28x28 pixels still retain key digit features of the image, what do you think?

thank you Jason

Jason Brownlee December 31, 2020 at 9:26 am #

REPLY ↗

Thanks.

Very nice experiments!

Yes, the larger images perhaps provide more context.

Hannes January 26, 2021 at 1:15 am #

Hi Jason,

You have always very interesting articles.

I was able to train a CNN based on the Char74K-dataset. It's about 42MB and the time to predict a number is fast.

I did use a kNN-model before (on the same dataset) though.

Thanks,

Hannes

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee January 26, 2021 at 5:58 am #

REPLY ↗

Thanks.

Well done!

Mukesh February 5, 2021 at 8:42 pm #

REPLY ↗

Hi Jason,

Great article!!

i just have one doubt, in the evaluate method why are we creating a new model for every fold?

Thanks,

Start Machine Learning

Jason Brownlee February 6, 2021 at 5:49 am #

REPLY ↗

This is how k-fold cross-validation works, we have a new training set and fit a new model for each fold, then average the performance of all models to get an estimate of model performance when used on unseen data.

You can learn more here:

<https://machinelearningmastery.com/k-fold-cross-validation/>

Dmitry February 6, 2021 at 6:20 pm #

REPLY ↗

Thanks for this amazing tutorial!

I have some experience in MLP but not with CNN. Is the data have shape 16×1 (row) with values between 0,01 for 15 classes. If I reshape samples from 16×1 to $4 \times 4 \times 1$ should I use? Or better dont reshape and try model with

Jason Brownlee February 7, 2021 at 5:17 am #

CNN can be effective if there is a spatial relationship between the inputs in an image or a sequence.

If you have tabular data (e.g. not images and not sequences), then a CNN does not make sense.

Ahmed Shafeek February 18, 2021 at 2:19 pm #

REPLY ↗

Hi Jason,

Thank you very much for your amazing tutorial. I am running your code using pycharm and I have a GPU with the right CUDA but your code is only running on my CPU.

Is there something that I am missing here?

thank you. 😊

Ahmed Shafeek February 18, 2021 at 4:52 pm #

REPLY ↗

Sorry Jason, it turned out I had a problem in the CuDNN installation and I did a work around to fix it and it worked. 😊

Thanks again for this great tutorial. 😊

[Start Machine Learning](#)

Jason Brownlee February 19, 2021 at 5:55 am #

REPLY ↗

No problem!

Happy to hear you fixed the issue.

Jason Brownlee February 19, 2021 at 5:54 am #

REPLY ↗

The code is agnostic to hardware – runs on both.

If it is running on your CPU, then you need to change the code to run on GPU. I have tutorials on this topic, sorry.

Karl February 24, 2021 at 1:20 am #

Hi Jason,

Thank you very much for this tutorial, helped a lot! I developed my own CNN, which seems to perform pretty well according to the mean accuracy. However, the cross-entropy loss of the validation set is below the loss of the training set. Is this a sign of underfitting? And if so, how can I fix it? I would like to test the model to unseen data?

Start Machine Learning

You can master applied Machine Learning without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee February 24, 2021 at 5:34 am #

REPLY ↗

You're welcome.

Well done!

If performance is poor, a learning curve can help diagnose issues, this can help you interpret the learning curve:

<https://machinelearningmastery.com/learning-curves-for-diagnosing-machine-learning-model-performance/>

Noah February 24, 2021 at 1:51 am #

REPLY ↗

Hi Jason,

Thank you very much for this tutorial, helped a lot! I developed my own CNN, which seems to perform pretty well according to the mean accuracy. However, the cross-entropy loss of the validation set is below the loss of the training set. Is this a sign of underfitting? And if so, how can I fix it? I would like to test the model to unseen data?

Start Machine Learning

Jason Brownlee February 24, 2021 at 5:35 am #

REPLY ↗

You're welcome.

Perhaps focus on the out of sample performance of the model first and optimize that.

Geo March 2, 2021 at 12:20 am #

REPLY ↗

Hi Jason,

Previously I was grateful for the tutorial you provided, I have one question though. Is it possible to get the probability value of the prediction result? (for example

Jason Brownlee March 2, 2021 at 5:45 am #

Thanks!

You can call `model.predict()` then multiply the result by 100 to get the percentage.

Start Machine Learning

You can master applied Machine Learning without **math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

William March 2, 2021 at 3:01 pm #

REPLY ↗

Thank you for the tutorial!

I created an image with same characteristics as `sample_image.png` (1490×1480, black background, white foreground), opened it in Paint, typed a “6” (without the quotes, of course), increased the font size to 1000 (so it would occupy most of the canvas), then saved it. However, the code doesn’t seem to recognize it. It seems to think it’s a “1”, “8”, or other numbers.

Any suggestions on why?

Jason Brownlee March 3, 2021 at 5:25 am #

REPLY ↗

Perhaps there was some important difference in the image itself or its preparation (e.g. pixel scaling) that differed from the training dataset?

Venkat April 1, 2021 at 10:07 pm #

REPLY ↗

Start Machine Learning

Hi Jason – fantastic article!! Just wondering why you had to do evaluate to get the accuracy, while it is already being returned by the fit call? I just checked the histories and they have exactly the same stuff returned by the evaluate call. Was there any other reason for the evaluate call?

Jason Brownlee April 2, 2021 at 5:39 am #

REPLY ↗

Thanks!

No, accuracy on a hold out dataset is not calculated when calling fit. We must make predictions on new data manually or use the evaluate function.

Venkat April 15, 2021 at 4:17 am #

Thanks for the reply!! I was seeing the evaluate calls for validation. As a result, the results returned by the evaluate call. Am I making sense?

```
# fit model
history = model.fit(trainX, trainY, epochs=10, b
verbose=0)
# evaluate model
_, acc = model.evaluate(testX, testY, verbose=0)
```

Start Machine Learning

X

You can master applied Machine Learning **without math or fancy degrees**.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Jason Brownlee April 15, 2021 at 5:30 am #

REPLY ↗

Generally it is not a good idea to use test data as validation, I do to keep the examples simple.

Berke April 4, 2021 at 9:47 pm #

REPLY ↗

Hi Jason, thanks for this great tutorial!

I don't get why we don't use the models that we created during the k-fold. I was expecting something like we combine this k models results into one single model and save it. If we don't do that how's that useful?

Jason Brownlee April 5, 2021 at 6:11 am #

REPLY ↗

Models created during k-fold cross-validation are discarded. They are only used to estimate the performance of the model/pipeline on unseen exar

Start Machine Learning

Once we choose a configuration, we can fit a final model on all data and use it to make predictions on new examples:

<https://machinelearningmastery.com/train-final-machine-learning-model/>

Vishwa April 15, 2021 at 9:38 pm #

REPLY ↗

Thanks a lot Jason !!! This tutorial was very very helpful and it solved half of my project problems !!!
Thanks a lot again !..

Jason Brownlee April 16, 2021 at 5:30 am #

You're welcome.

ismail June 9, 2021 at 8:47 am #

Thank you for the great tutorial. Is it possible to

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

REPLY ↗

Jason Brownlee June 10, 2021 at 5:20 am #

Perhaps try it and see.

Taylor Rayne June 12, 2021 at 12:20 pm #

REPLY ↗

Hello! Thank you Jason for this amazing resource. I am currently working on a rented computer from school – do you think I could do implement this code in CoLabs?

Jason Brownlee June 13, 2021 at 5:46 am #

REPLY ↗

You're welcome.

Good question, see this:

<https://machinelearningmastery.com/faq/single-faq/do-code-examples-run-on-google-colab>

Manal June 15, 2021 at 4:14 pm #

Start Machine Learning

Hi Jason,
Thank you very much
If possible, advise on the implementation HMM.

Jason Brownlee June 16, 2021 at 6:17 am #

REPLY ↗

Thanks for the suggestion, perhaps in the future.

amberbir June 15, 2021 at 11:11 pm #

REPLY ↗

i am developing a model of handwritten chara
is not only because i have a small datatset, there is so
improvement

Jason Brownlee June 16, 2021 at 6:21 am #

Perhaps try some of these suggestions:
<https://machinelearningmastery.com/improve-deep/>

Start Machine Learning

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Shri July 8, 2021 at 11:15 pm #

REPLY ↗

Can I use this code for any character dataset ?

Jason Brownlee July 9, 2021 at 5:11 am #

REPLY ↗

Perhaps try it and see.

Shireesh Apte August 8, 2021 at 8:34 am #

REPLY ↗

Could you provide a peer review for a manuscript that has been submitted to the Journal of High School Science (<http://jhss.scholasticahq.com>) ? The title of the manuscript is : Implementing a Quantum Convolutional Neural Network for Efficient Image Recognition. Abstract: Machine learning has many real-world applications ranging from modeling the universe to computational chemistry. As probability is the bedrock for machine learning, it is essential to optimize both hardware and software to obtain the best results. Classical computers are generally used for machine learning programs. However, learning from high-dimensional data often demands excessive comp

Start Machine Learning

highest accuracy. The Quantum Computing environment can be utilized to create a more accurate model than that created via classical computing. To test this quantum advantage, we implemented a Quantum Convolutional Neural Network (QCNN), which parallels the structure of the classical Convolution Neural Network (CNN) in the quantum domain. Due to the lack of quantum computers with many qubits, physicists, namely John Preskill, have introduced the Noisy Intermediate Scale Quantum (NISQ) concept, which constitutes a hybrid interface between classical and quantum computers. In the context of QCNN, the data processing and the cost function optimization would be performed on the classical computer, while the probabilities generated by the Variational Quantum Circuits (VQC) would be evaluated on the quantum computer. The QCNN consists of a classical-to-quantum data encoder, a cluster state quantum circuit to entangle qubit states, a series of Variational Quantum Circuits using Quantum Convolutional and Pooling Layers for efficient feature extraction, a quantum-to-classical data decoder, which would lead to the output. Both the CNN and the QCNN extract features from data like 2D images, and performances can be compared using metrics like accuracy and time.

Jason Brownlee August 9, 2021 at 5:52 am #

Sorry I cannot peer review your manuscri

Leave a Reply

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

Email Address

START MY EMAIL COURSE

Name (required)

Email (will not be published) (required)

Website

SUBMIT COMMENT

Start Machine Learning

Welcome!

I'm Jason Brownlee PhD
and I **help developers** get results with **machine learning**.
[Read more](#)

Never miss a tutorial:**Picked for you:**

[How to Train an Object Detection Model with Keras](#)

[How to Develop a Face Recognition System Using FaceNet](#)

[How to Perform Object Detection With YOLOv3 in Keras](#)

[How to Classify Photos of Dogs and Cats \(with 97% accuracy\)](#)

[How to Get Started With Deep Learning for Computer Vision](#)

Start Machine Learning X

You can master applied Machine Learning
without math or fancy degrees.

Find out how in this *free* and *practical* course.

[START MY EMAIL COURSE](#)

Loving the Tutorials?

The Deep Learning for Computer Vision EBook is where you'll find the **Really Good** stuff.

[>> SEE WHAT'S INSIDE](#)

© 2021 Machine Learning Mastery Pty. Ltd. All Rights Reserved.

[LinkedIn](#) | [Twitter](#) | [Facebook](#) | [Newsletter](#) | [RSS](#)

[Privacy](#) | [Disclaimer](#) | [Terms](#) | [Contact](#) | [Sitemap](#) | [Search](#)

[Start Machine Learning](#)