

## ▼ Transfer learning / fine-tuning

This tutorial will guide you through the process of using *transfer learning* to learn an accurate image classifier from a relatively small number of training samples. Generally speaking, transfer learning refers to the process of leveraging the knowledge learned in one model for the training of another model.

More specifically, the process involves taking an existing neural network which was previously trained to good performance on a larger dataset, and using it as the basis for a new model which leverages that previous network's accuracy for a new task. This method has become popular in recent years to improve the performance of a neural net trained on a small dataset; the intuition is that the new dataset may be too small to train to good performance by itself, but we know that most neural nets trained to learn image features often learn similar features anyway, especially at early layers where they are more generic (edge detectors, blobs, and so on).

Transfer learning has been largely enabled by the open-sourcing of state-of-the-art models; for the top performing models in image classification tasks (like from [ILSVRC](#)), it is common practice now to not only publish the architecture, but to release the trained weights of the model as well. This lets amateurs use these top image classifiers to boost the performance of their own task-specific models.

### Feature extraction vs. fine-tuning

At one extreme, transfer learning can involve taking the pre-trained network and freezing the weights, and using one of its hidden layers (usually the last one) as a feature extractor, using those features as the input to a smaller neural net.

At the other extreme, we start with the pre-trained network, but we allow some of the weights (usually the last layer or last few layers) to be modified. Another name for this procedure is called "fine-tuning" because we are slightly adjusting the pre-trained net's weights to the new task. We usually train such a network with a lower learning rate, since we expect the features are already relatively good and do not need to be changed too much.

Sometimes, we do something in-between: Freeze just the early/generic layers, but fine-tune the later layers. Which strategy is best depends on the size of your dataset, the number of classes, and how much it resembles the dataset the previous model was trained on (and thus, whether it can benefit from the same learned feature extractors). A more detailed discussion of how to strategize can be found in [1] [2].

## Procedure

In this guide will go through the process of loading a state-of-the-art, 1000-class image classifier, [VGG16](#) which [won the ImageNet challenge in 2014](#), and using it as a fixed feature

extractor to train a smaller custom classifier on our own images, although with very few code changes, you can try fine-tuning as well.

We will first load VGG16 and remove its final layer, the 1000-class softmax classification layer specific to ImageNet, and replace it with a new classification layer for the classes we are training over. We will then freeze all the weights in the network except the new ones connecting to the new classification layer, and then train the new classification layer over our new dataset.

We will also compare this method to training a small neural network from scratch on the new dataset, and as we shall see, it will dramatically improve our accuracy. We will do that part first.

As our test subject, we'll use a dataset consisting of around 6000 images belonging to 97 classes, and train an image classifier with around 80% accuracy on it. It's worth noting that this strategy scales well to image sets where you may have even just a couple hundred or less images. Its performance will be lesser from a small number of samples (depending on classes) as usual, but still impressive considering the usual constraints.

```
%matplotlib inline

import os

#if using Theano with GPU
#os.environ["KERAS_BACKEND"] = "tensorflow"

import random
import numpy as np
import keras

import matplotlib.pyplot as plt
from matplotlib.pyplot import imshow

from keras.preprocessing import image
from keras.applications.imagenet_utils import preprocess_input
from keras.models import Sequential
from keras.layers import Dense, Dropout, Flatten, Activation
from keras.layers import Conv2D, MaxPooling2D
from keras.models import Model

Using TensorFlow backend.
```

## ▼ Getting a dataset

The first step is going to be to load our data. As our example, we will be using the dataset [CalTech-101](#), which contains around 9000 labeled images belonging to 101 object categories. However, we will exclude 5 of the categories which have the most images. This is in order to keep the class distribution fairly balanced (around 50-100) and constrained to a smaller number of images, around 6000.

To obtain this dataset, you can either run the download script `download.sh` in the `data` folder, or the following commands:

```
wget http://www.vision.caltech.edu/Image_Datasets/Caltech101/101_ObjectCategories.tar.gz  
tar -xvzf 101_ObjectCategories.tar.gz
```

If you wish to use your own dataset, it should be arranged in the same fashion to 101\_ObjectCategories with all of the images organized into subfolders, one for each class. In this case, the following cell should load your custom dataset correctly by just replacing root with your folder. If you have an alternate structure, you just need to make sure that you load the list data where every element is a dict where x is the data (a 1-d numpy array) and y is the label (an integer). Use the helper function get\_image(path) to load the image correctly into the array, and note also that the images are being resized to 224x224. This is necessary because the input to VGG16 is a 224x224 RGB image. You do not need to resize them on your hard drive, as that is being done in the code below.

If you have 101 ObjectCategories in your data folder, the following cell should load all the data.

```
!echo "Downloading 101_Object_Categories for image notebooks"
!curl -L -o 101_ObjectCategories.tar.gz --progress-bar http://www.vision.caltech.edu/Image
!tar -xzf 101_ObjectCategories.tar.gz
!rm 101_ObjectCategories.tar.gz
!ls
```

Downloading 101\_Object\_Categories for image notebooks  
#####
101 ObjectCategories sample data 100.0%

```
root = '101_ObjectCategories'
exclude = ['BACKGROUND_Google', 'Motorbikes', 'airplanes', 'Faces_easy', 'Faces']
train_split, val_split = 0.7, 0.15

categories = [x[0] for x in os.walk(root) if x[0]][1:]
categories = [c for c in categories if c not in [os.path.join(root, e) for e in exclude]]

print(categories)
```

This function is useful for pre-processing the data into an image and input vector.

```
# helper function to load image and return it and input vector
def get_image(path):
    img = image.load_img(path, target_size=(224, 224))
    x = image.img_to_array(img)
```

```
x = np.expand_dims(x, axis=0)
x = preprocess_input(x)
return img, x
```

Load all the images from root folder

```
data = []
for c, category in enumerate(categories):
    images = [os.path.join(dp, f) for dp, dn, filenames
              in os.walk(category) for f in filenames
              if os.path.splitext(f)[1].lower() in ['.jpg', '.png', '.jpeg']]
    for img_path in images:
        img, x = get_image(img_path)
        data.append({'x':np.array(x[0]), 'y':c})

# count the number of classes
num_classes = len(categories)
```

Randomize the data order.

```
random.shuffle(data)
```

create training / validation / test split (70%, 15%, 15%)

```
idx_val = int(train_split * len(data))
idx_test = int((train_split + val_split) * len(data))
train = data[:idx_val]
val = data[idx_val:idx_test]
test = data[idx_test:]
```

Separate data for labels.

```
x_train, y_train = np.array([t["x"] for t in train]), [t["y"] for t in train]
x_val, y_val = np.array([t["x"] for t in val]), [t["y"] for t in val]
x_test, y_test = np.array([t["x"] for t in test]), [t["y"] for t in test]
print(y_test)
```

```
[51, 24, 43, 74, 5, 14, 33, 57, 40, 14, 63, 50, 28, 91, 28, 43, 53, 28, 81, 93, 25, 6]
```

Pre-process the data as before by making sure it's float32 and normalized between 0 and 1.

```
# normalize data
x_train = x_train.astype('float32') / 255.
x_val = x_val.astype('float32') / 255.
x_test = x_test.astype('float32') / 255.
```

```
# convert labels to one-hot vectors
y_train = keras.utils.to_categorical(y_train, num_classes)
y_val = keras.utils.to_categorical(y_val, num_classes)
y_test = keras.utils.to_categorical(y_test, num_classes)
print(y_test.shape)

(932, 97)
```

Let's get a summary of what we have.

```
# summary
print("finished loading %d images from %d categories"%(len(data), num_classes))
print("train / validation / test split: %d, %d, %d"%(len(x_train), len(x_val), len(x_test)))
print("training data shape: ", x_train.shape)
print("training labels shape: ", y_train.shape)

finished loading 6209 images from 97 categories
train / validation / test split: 4346, 931, 932
training data shape:  (4346, 224, 224, 3)
training labels shape:  (4346, 97)
```

If everything worked properly, you should have loaded a bunch of images, and split them into three sets: `train`, `val`, and `test`. The shape of the training data should be  $(n, 224, 224, 3)$  where  $n$  is the size of your training set, and the labels should be  $(n, c)$  where  $c$  is the number of classes (97 in the case of `101_ObjectCategories`).

Notice that we divided all the data into three subsets – a training set `train`, a validation set `val`, and a test set `test`. The reason for this is to properly evaluate the accuracy of our classifier. During training, the optimizer uses the validation set to evaluate its internal performance, in order to determine the gradient without overfitting to the training set. The `test` set is always held out from the training algorithm, and is only used at the end to evaluate the final accuracy of our model.

Let's quickly look at a few sample images from our dataset.

```
images = [os.path.join(dp, f) for dp, dn, filenames in os.walk(root) for f in filenames if
idx = [int(len(images) * random.random()) for i in range(8)]
imgs = [image.load_img(images[i], target_size=(224, 224)) for i in idx]
concat_image = np.concatenate([np.asarray(img) for img in imgs], axis=1)
plt.figure(figsize=(16,4))
plt.imshow(concat_image)
```

```
<matplotlib.image.AxesImage at 0x7f233bce49b0>
```



## ▼ First training a neural net from scratch

Before doing the transfer learning, let's first build a neural network from scratch for doing classification on our dataset. This will give us a baseline to compare to our transfer-learned network later.

The network we will construct contains 4 alternating convolutional and max-pooling layers, followed by a dropout after every other conv/pooling pair. After the last pooling layer, we will attach a fully-connected layer with 256 neurons, another dropout layer, then finally a softmax classification layer for our classes.

Our loss function will be, as usual, categorical cross-entropy loss, and our learning algorithm will be AdaDelta. Various things about this network can be changed to get better performance, perhaps using a larger network or a different optimizer will help, but for the purposes of this notebook, the goal is to just get an understanding of an approximate baseline for comparison's sake, and so it isn't necessary to spend much time trying to optimize this network.

Upon compiling the network, let's run `model.summary()` to get a snapshot of its layers.

```
#·build·the·network
model=·Sequential()
print("Input·dimensions:·",x_train.shape[1:])

model.add(Conv2D(32,·(3,·3),·input_shape=x_train.shape[1:]))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2,·2)))

model.add(Conv2D(32,·(3,·3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2,·2)))

model.add(Dropout(0.25))

model.add(Conv2D(32,·(3,·3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2,·2)))

model.add(Conv2D(32,·(3,·3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2,·2)))

model.add(Dropout(0.25))

model.add(Flatten())
model.add(Dense(256))
model.add(Activation('relu'))

model.add(Dropout(0.5))
```

```
model.add(Dense(num_classes))
model.add(Activation('softmax'))
```

```
model.summary()
```

Input dimensions: (224, 224, 3)

Layer (type)	Output Shape	Param #
<hr/>		
conv2d_17 (Conv2D)	(None, 222, 222, 32)	896
activation_25 (Activation)	(None, 222, 222, 32)	0
max_pooling2d_17 (MaxPooling)	(None, 111, 111, 32)	0
conv2d_18 (Conv2D)	(None, 109, 109, 32)	9248
activation_26 (Activation)	(None, 109, 109, 32)	0
max_pooling2d_18 (MaxPooling)	(None, 54, 54, 32)	0
dropout_13 (Dropout)	(None, 54, 54, 32)	0
conv2d_19 (Conv2D)	(None, 52, 52, 32)	9248
activation_27 (Activation)	(None, 52, 52, 32)	0
max_pooling2d_19 (MaxPooling)	(None, 26, 26, 32)	0
conv2d_20 (Conv2D)	(None, 24, 24, 32)	9248
activation_28 (Activation)	(None, 24, 24, 32)	0
max_pooling2d_20 (MaxPooling)	(None, 12, 12, 32)	0
dropout_14 (Dropout)	(None, 12, 12, 32)	0
flatten_5 (Flatten)	(None, 4608)	0
dense_9 (Dense)	(None, 256)	1179904
activation_29 (Activation)	(None, 256)	0
dropout_15 (Dropout)	(None, 256)	0
dense_10 (Dense)	(None, 97)	24929
activation_30 (Activation)	(None, 97)	0
<hr/>		
Total params:	1,233,473	
Trainable params:	1,233,473	
Non-trainable params:	0	

We've created a medium-sized network with ~1.2 million weights and biases (the parameters). Most of them are leading into the one pre-softmax fully-connected layer "dense\_5".

We can now go ahead and train our model for 100 epochs with a batch size of 128. We'll also record its history so we can plot the loss over time later.

```
# compile the model to use categorical cross-entropy loss function and adadelta optimizer
model.compile(loss='categorical_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])

history = model.fit(x_train, y_train,
                     batch_size=128,
                     epochs=10,
                     validation_data=(x_val, y_val))

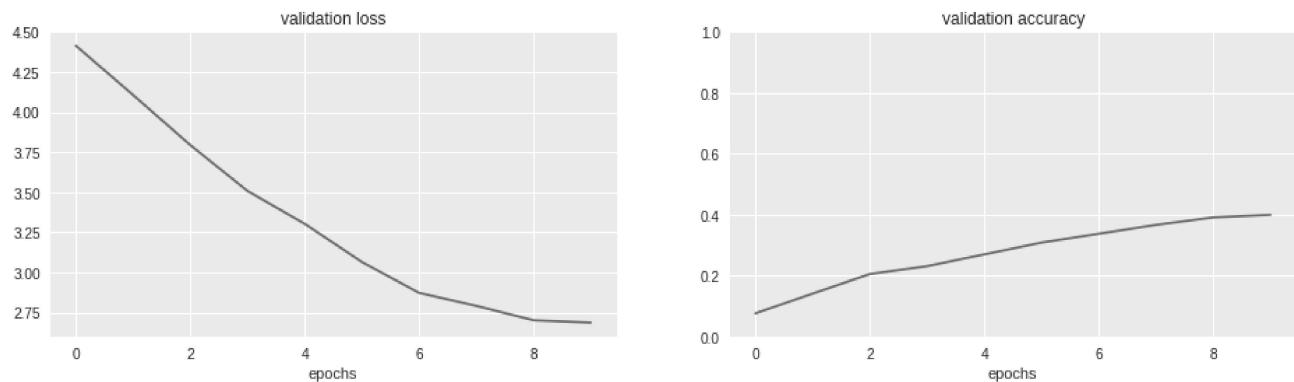
Train on 4346 samples, validate on 931 samples
Epoch 1/10
4346/4346 [=====] - 11s 3ms/step - loss: 4.5153 - acc: 0.046
Epoch 2/10
4346/4346 [=====] - 10s 2ms/step - loss: 4.2447 - acc: 0.101
Epoch 3/10
4346/4346 [=====] - 10s 2ms/step - loss: 3.9434 - acc: 0.151
Epoch 4/10
4346/4346 [=====] - 10s 2ms/step - loss: 3.6178 - acc: 0.204
Epoch 5/10
4346/4346 [=====] - 10s 2ms/step - loss: 3.3217 - acc: 0.251
Epoch 6/10
4346/4346 [=====] - 10s 2ms/step - loss: 3.0783 - acc: 0.291
Epoch 7/10
4346/4346 [=====] - 10s 2ms/step - loss: 2.7861 - acc: 0.356
Epoch 8/10
4346/4346 [=====] - 10s 2ms/step - loss: 2.5162 - acc: 0.391
Epoch 9/10
4346/4346 [=====] - 10s 2ms/step - loss: 2.2286 - acc: 0.441
Epoch 10/10
4346/4346 [=====] - 10s 2ms/step - loss: 2.0241 - acc: 0.491
```

Let's plot the validation loss and validation accuracy over time.

```
fig = plt.figure(figsize=(16,4))
ax = fig.add_subplot(121)
ax.plot(history.history["val_loss"])
ax.set_title("validation loss")
ax.set_xlabel("epochs")

ax2 = fig.add_subplot(122)
ax2.plot(history.history["val_acc"])
ax2.set_title("validation accuracy")
ax2.set_xlabel("epochs")
ax2.set_ylim(0, 1)

plt.show()
```



Notice that the validation loss begins to actually rise after around 16 epochs, even though validation accuracy remains roughly between 40% and 50%. This suggests our model begins overfitting around then, and best performance would have been achieved if we had stopped early around then. Nevertheless, our accuracy would not have likely been above 50%, and probably lower down.

We can also get a final evaluation by running our model on the training set. Doing so, we get the following results:

```
loss, accuracy = model.evaluate(x_test, y_test, verbose=0)
print('Test loss:', loss)
print('Test accuracy:', accuracy)

Test loss: 2.6905115888865723
Test accuracy: 0.4034334763948498
```

Finally, we see that we have achieved a (top-1) accuracy of around 49%. That's not too bad for 6000 images, considering that if we were to use a naive strategy of taking random guesses, we would have only gotten around 1% accuracy.

## ▼ Transfer learning by starting with existing network

Now we can move on to the main strategy for training an image classifier on our small dataset: by starting with a larger and already trained network.

To start, we will load the VGG16 from keras, which was trained on ImageNet and the weights saved online. If this is your first time loading VGG16, you'll need to wait a bit for the weights to download from the web. Once the network is loaded, we can again inspect the layers with the `summary()` method.

```
vgg = keras.applications.VGG16(weights='imagenet', include_top=True)
vgg.summary()
```

Downloading data from <https://github.com/fchollet/deep-learning-models/releases/download/v553467904/553467096> [=====] - 7s 0us/step

Layer (type)	Output Shape	Param #
<hr/>		
input_1 (InputLayer)	(None, 224, 224, 3)	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312
predictions (Dense)	(None, 1000)	4097000
<hr/>		

Total params: 138,357,544

Trainable params: 138,357,544

Non-trainable params: 0

Notice that VGG16 is *much* bigger than the network we constructed earlier. It contains 13 convolutional layers and two fully connected layers at the end, and has over 138 million parameters, around 100 times as many parameters than the network we made above. Like our first network, the majority of the parameters are stored in the connections leading into the first fully-connected layer.

VGG16 was made to solve ImageNet, and achieves a 8.8% top-5 error rate, which means that 91.2% of test samples were classified correctly within the top 5 predictions for each image. Its top-1 accuracy--equivalent to the accuracy metric we've been using (that the top prediction is correct)--is 73%. This is especially impressive since there are not just 97, but 1000 classes, meaning that random guesses would get us only 0.1% accuracy.

In order to use this network for our task, we "remove" the final classification layer, the 1000-neuron softmax layer at the end, which corresponds to ImageNet, and instead replace it with a new softmax layer for our dataset, which contains 97 neurons in the case of the 101\_ObjectCategories dataset.

In terms of implementation, it's easier to simply create a copy of VGG from its input layer until the second to last layer, and then work with that, rather than modifying the VGG object directly. So technically we never "remove" anything, we just circumvent/ignore it. This can be done in the following way, by using the keras `Model` class to initialize a new model whose input layer is the same as VGG but whose output layer is our new softmax layer, called `new_classification_layer`. Note: although it appears we are duplicating this large network, internally Keras is actually just copying all the layers by reference, and thus we don't need to worry about overloading the memory.

```
# make a reference to VGG's input layer
inp = vgg.input

# make a new softmax layer with num_classes neurons
new_classification_layer = Dense(num_classes, activation='softmax')

# connect our new layer to the second to last layer in VGG, and make a reference to it
out = new_classification_layer(vgg.layers[-2].output)

# create a new network between inp and out
model_new = Model(inp, out)
```

We are going to retrain this network, `model_new` on the new dataset and labels. But first, we need to freeze the weights and biases in all the layers in the network, except our new one at the end, with the expectation that the features that were learned in VGG should still be fairly relevant to the new image classification task. Not optimal, but most likely better than what we can train to in our limited dataset.

By setting the trainable flag in each layer false (except our new classification layer), we ensure all the weights and biases in those layers remain fixed, and we simply train the weights in the one layer at the end. In some cases, it is desirable to *not* freeze all the pre-classification layers. If your dataset has enough samples, and doesn't resemble ImageNet very much, it might be advantageous to fine-tune some of the VGG layers along with the new classifier, or possibly even all of them. To do this, you can change the below code to make more of the layers trainable.

In the case of CalTech-101, we will just do feature extraction, fearing that fine-tuning too much with this dataset may overfit. But maybe we are wrong? A good exercise would be to try out both, and compare the results.

So we go ahead and freeze the layers, and compile the new model with exactly the same optimizer and loss function as in our first network, for the sake of a fair comparison. We then run summary again to look at the network's architecture.

```
# make all layers untrainable by freezing weights (except for last layer)
for l, layer in enumerate(model_new.layers[:-1]):
    layer.trainable = False

# ensure the last layer is trainable/not frozen
for l, layer in enumerate(model_new.layers[-1:]):
    layer.trainable = True

model_new.compile(loss='categorical_crossentropy',
                  optimizer='adam',
                  metrics=['accuracy'])

model_new.summary()
```

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	(None, 224, 224, 3)	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0

block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312
dense_12 (Dense)	(None, 97)	397409
<hr/>		
Total params: 134,657,953		
Trainable params: 397,409		
Non-trainable params: 134,260,544		

Looking at the summary, we see the network is identical to the VGG model we instantiated earlier, except the last layer, formerly a 1000-neuron softmax, has been replaced by a new 97-neuron softmax. Additionally, we still have roughly 134 million weights, but now the vast majority of them are "non-trainable params" because we froze the layers they are contained in. We now only have 397,000 trainable parameters, which is actually only a quarter of the number of parameters needed to train the first model.

As before, we go ahead and train the new model, using the same hyperparameters (batch size and number of epochs) as before, along with the same optimization algorithm. We also keep track of its history as we go.

```
history2 = model_new.fit(x_train, y_train,
                          batch_size=128,
                          epochs=10,
                          validation_data=(x_val, y_val))
```

```
Train on 4346 samples, validate on 931 samples
Epoch 1/10
4346/4346 [=====] - 66s 15ms/step - loss: 4.0574 - acc: 0.16
Epoch 2/10
4346/4346 [=====] - 46s 11ms/step - loss: 2.5812 - acc: 0.44
Epoch 3/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.9747 - acc: 0.56
```

```
Epoch 4/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.6366 - acc: 0.61
Epoch 5/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.4263 - acc: 0.61
Epoch 6/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.2610 - acc: 0.76
Epoch 7/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.1283 - acc: 0.72
Epoch 8/10
4346/4346 [=====] - 46s 11ms/step - loss: 1.0551 - acc: 0.71
Epoch 9/10
4346/4346 [=====] - 46s 11ms/step - loss: 0.9447 - acc: 0.78
Epoch 10/10
4346/4346 [=====] - 46s 11ms/step - loss: 0.8976 - acc: 0.79
```



Our validation accuracy hovers close to 80% towards the end, which is more than 30% improvement on the original network trained from scratch (meaning that we make the wrong prediction on 20% of samples, rather than 50%).

It's worth noting also that this network actually trains *slightly faster* than the original network, despite having more than 100 times as many parameters! This is because freezing the weights negates the need to backpropagate through all those layers, saving us on runtime.

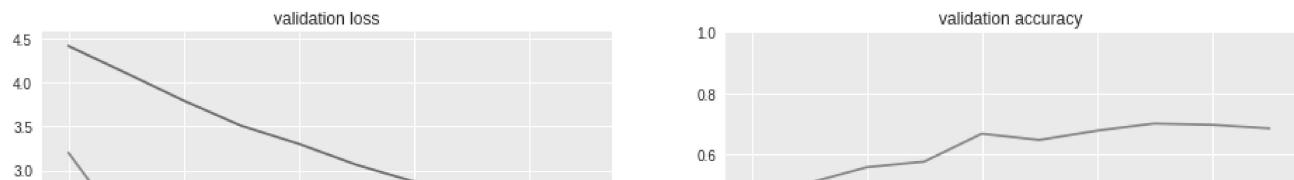
Let's plot the validation loss and accuracy again, this time comparing the original model trained from scratch (in blue) and the new transfer-learned model in green.

```
fig = plt.figure(figsize=(16,4))
ax = fig.add_subplot(121)
ax.plot(history.history["val_loss"])
ax.plot(history2.history["val_loss"])
ax.set_title("validation loss")
ax.set_xlabel("epochs")

ax2 = fig.add_subplot(122)
ax2.plot(history.history["val_acc"])
ax2.plot(history2.history["val_acc"])
ax2.set_title("validation accuracy")
ax2.set_xlabel("epochs")
ax2.set_xlim(0, 1)

plt.show()
```





Notice that whereas the original model began overfitting around epoch 16, the new model continued to slowly decrease its loss over time, and likely would have improved its accuracy slightly with more iterations. The new model made it to roughly 80% top-1 accuracy (in the validation set) and continued to improve slowly through 100 epochs.

It's possibly we could have improved the original model with better regularization or more dropout, but we surely would not have made up the >30% improvement in accuracy.

Again, we do a final validation on the test set.

```
loss, ·accuracy··=·model_new.evaluate(x_test, ·y_test, ·verbose=0)
```

```
print('Test·loss:', ·loss)
print('Test·accuracy:', ·accuracy)
```

```
Test loss: 1.1540323304004423
Test accuracy: 0.7156652360515021
```

To predict a new image, simply run the following code to get the probabilities for each class.

```
img, x = get_image('101_ObjectCategories/airplanes/image_0003.jpg')
probabilities = model_new.predict([x])
```

```
0.6425913
```

## Improving the results

78.2% top-1 accuracy on 97 classes, roughly evenly distributed, is a pretty good achievement. It is not quite as impressive as the original VGG16 which achieved 73% top-1 accuracy on 1000 classes. Nevertheless, it is much better than what we were able to achieve with our original network, and there is room for improvement. Some techniques which possibly could have improved our performance.

- Using data augmentation: augmentation refers to using various modifications of the original training data, in the form of distortions, rotations, rescalings, lighting changes, etc to increase the size of the training set and create more tolerance for such distortions.
- Using a different optimizer, adding more regularization/dropout, and other hyperparameters.
- Training for longer (of course)

A more advanced example of transfer learning in Keras, involving augmentation for a small 2-class dataset, can be found in the [Keras blog](#).

Produtos pagos do Colab - Cancelar contratos

