

✓ Transfer learning & fine-tuning

Author: [fchollet](#)

Date created: 2020/04/15

Last modified: 2023/06/25

Description: Complete guide to transfer learning & fine-tuning in Keras.

✓ Setup

```
1 import numpy as np
2 import keras
3 from keras import layers
4 import tensorflow_datasets as tfds
5 import matplotlib.pyplot as plt
```

Introduction

Transfer learning consists of taking features learned on one problem, and leveraging them on a new, similar problem. For instance, features from a model that has learned to identify racoons may be useful to kick-start a model meant to identify tanukis.

Transfer learning is usually done for tasks where your dataset has too little data to train a full-scale model from scratch.

The most common incarnation of transfer learning in the context of deep learning is the following workflow:

1. Take layers from a previously trained model.
2. Freeze them, so as to avoid destroying any of the information they contain during future training rounds.
3. Add some new, trainable layers on top of the frozen layers. They will learn to turn the old features into predictions on a new dataset.
4. Train the new layers on your dataset.

A last, optional step, is **fine-tuning**, which consists of unfreezing the entire model you obtained above (or part of it), and re-training it on the new data with a very low learning rate. This can potentially achieve meaningful improvements, by incrementally adapting the pretrained

features to the new data.

First, we will go over the Keras `trainable` API in detail, which underlies most transfer learning & fine-tuning workflows.

Then, we'll demonstrate the typical workflow by taking a model pretrained on the ImageNet dataset, and retraining it on the Kaggle "cats vs dogs" classification dataset.

This is adapted from [Deep Learning with Python](#) and the 2016 blog post ["building powerful image classification models using very little data"](#).

✓ Freezing layers: understanding the `trainable` attribute

Layers & models have three weight attributes:

- `weights` is the list of all weights variables of the layer.
- `trainable_weights` is the list of those that are meant to be updated (via gradient descent) to minimize the loss during training.
- `non_trainable_weights` is the list of those that aren't meant to be trained. Typically they are updated by the model during the forward pass.

Example: the `Dense` layer has 2 trainable weights (kernel & bias)

```
1 layer = keras.layers.Dense(3)
2 layer.build((None, 4)) # Create the weights
3
4 print("weights:", len(layer.weights))
5 print("trainable_weights:", len(layer.trainable_weights))
6 print("non_trainable_weights:", len(layer.non_trainable_weights))
```

In general, all weights are trainable weights. The only built-in layer that has non-trainable weights is the `BatchNormalization` layer. It uses non-trainable weights to keep track of the mean and variance of its inputs during training. To learn how to use non-trainable weights in your own custom layers, see the [guide to writing new layers from scratch](#).

Example: the `BatchNormalization` layer has 2 trainable weights and 2 non-trainable weights

```
1 layer = keras.layers.BatchNormalization()
2 layer.build((None, 4)) # Create the weights
3
4 print("weights:", len(layer.weights))
5 print("trainable_weights:", len(layer.trainable_weights))
6 print("non_trainable_weights:", len(layer.non_trainable_weights))
```

Layers & models also feature a boolean attribute `trainable`. Its value can be changed. Setting `layer.trainable` to `False` moves all the layer's weights from trainable to non-trainable. This is called "freezing" the layer: the state of a frozen layer won't be updated during training (either when training with `fit()` or when training with any custom loop that relies on `trainable_weights` to apply gradient updates).

Example: setting trainable to False

```
1 layer = keras.layers.Dense(3)
2 layer.build((None, 4)) # Create the weights
3 layer.trainable = False # Freeze the layer
4
5 print("weights:", len(layer.weights))
6 print("trainable_weights:", len(layer.trainable_weights))
7 print("non_trainable_weights:", len(layer.non_trainable_weights))
```

When a trainable weight becomes non-trainable, its value is no longer updated during training.

```
1 # Make a model with 2 layers
2 layer1 = keras.layers.Dense(3, activation="relu")
3 layer2 = keras.layers.Dense(3, activation="sigmoid")
4 model = keras.Sequential([keras.Input(shape=(3,)), layer1, layer2])
5
6 # Freeze the first layer
7 layer1.trainable = False
8
9 # Keep a copy of the weights of layer1 for later reference
10 initial_layer1_weights_values = layer1.get_weights()
11
12 # Train the model
13 model.compile(optimizer="adam", loss="mse")
14 model.fit(np.random.random((2, 3)), np.random.random((2, 3)))
15
16 # Check that the weights of layer1 have not changed during training
17 final_layer1_weights_values = layer1.get_weights()
18 np.testing.assert_allclose(
19     initial_layer1_weights_values[0], final_layer1_weights_values[0]
20 )
21 np.testing.assert_allclose(
22     initial_layer1_weights_values[1], final_layer1_weights_values[1]
23 )
```

Do not confuse the `layer.trainable` attribute with the argument `training` in `layer.__call__()` (which controls whether the layer should run its forward pass in inference mode or training mode). For more information, see the [Keras FAQ](#).

✓ Recursive setting of the trainable attribute

If you set `trainable = False` on a model or on any layer that has sublayers, all children layers become non-trainable as well.

Example:

```
1 inner_model = keras.Sequential(  
2     [  
3         keras.Input(shape=(3,)),  
4         keras.layers.Dense(3, activation="relu"),  
5         keras.layers.Dense(3, activation="relu"),  
6     ]  
7 )  
8  
9 model = keras.Sequential(  
10     [  
11         keras.Input(shape=(3,)),  
12         inner_model,  
13         keras.layers.Dense(3, activation="sigmoid"),  
14     ]  
15 )  
16  
17 model.trainable = False # Freeze the outer model  
18  
19 assert inner_model.trainable == False # All layers in `model` are now frozen  
20 assert inner_model.layers[0].trainable == False # `trainable` is propagated recursively
```

The typical transfer-learning workflow

This leads us to how a typical transfer learning workflow can be implemented in Keras:

1. Instantiate a base model and load pre-trained weights into it.
2. Freeze all layers in the base model by setting `trainable = False`.
3. Create a new model on top of the output of one (or several) layers from the base model.
4. Train your new model on your new dataset.

Note that an alternative, more lightweight workflow could also be:

1. Instantiate a base model and load pre-trained weights into it.
2. Run your new dataset through it and record the output of one (or several) layers from the base model. This is called **feature extraction**.
3. Use that output as input data for a new, smaller model.

A key advantage of that second workflow is that you only run the base model once on your data, rather than once per epoch of training. So it's a lot faster & cheaper.

An issue with that second workflow, though, is that it doesn't allow you to dynamically modify the input data of your new model during training, which is required when doing data augmentation, for instance. Transfer learning is typically used for tasks when your new dataset has too little data to train a full-scale model from scratch, and in such scenarios data augmentation is very important. So in what follows, we will focus on the first workflow.

Here's what the first workflow looks like in Keras:

First, instantiate a base model with pre-trained weights.

```
base_model = keras.applications.Xception(  
    weights='imagenet', # Load weights pre-trained on ImageNet.  
    input_shape=(150, 150, 3),  
    include_top=False) # Do not include the ImageNet classifier at the top.
```

Then, freeze the base model.

```
base_model.trainable = False
```

Create a new model on top.

```
inputs = keras.Input(shape=(150, 150, 3))  
# We make sure that the base_model is running in inference mode here,  
# by passing `training=False`. This is important for fine-tuning, as you will  
# learn in a few paragraphs.  
x = base_model(inputs, training=False)  
# Convert features of shape `base_model.output_shape[1:]` to vectors  
x = keras.layers.GlobalAveragePooling2D()(x)  
# A Dense classifier with a single unit (binary classification)  
outputs = keras.layers.Dense(1)(x)  
model = keras.Model(inputs, outputs)
```

Train the model on new data.

```
model.compile(optimizer=keras.optimizers.Adam(),
              loss=keras.losses.BinaryCrossentropy(from_logits=True),
              metrics=[keras.metrics.BinaryAccuracy()])
model.fit(new_dataset, epochs=20, callbacks=..., validation_data=...)
```

Fine-tuning

Once your model has converged on the new data, you can try to unfreeze all or part of the base model and retrain the whole model end-to-end with a very low learning rate.

This is an optional last step that can potentially give you incremental improvements. It could also potentially lead to quick overfitting -- keep that in mind.

It is critical to only do this step *after* the model with frozen layers has been trained to convergence. If you mix randomly-initialized trainable layers with trainable layers that hold pre-trained features, the randomly-initialized layers will cause very large gradient updates during training, which will destroy your pre-trained features.

It's also critical to use a very low learning rate at this stage, because you are training a much larger model than in the first round of training, on a dataset that is typically very small. As a result, you are at risk of overfitting very quickly if you apply large weight updates. Here, you only want to readapt the pretrained weights in an incremental way.

This is how to implement fine-tuning of the whole base model:

```
# Unfreeze the base model
base_model.trainable = True

# It's important to recompile your model after you make any changes
# to the `trainable` attribute of any inner layer, so that your changes
# are take into account
model.compile(optimizer=keras.optimizers.Adam(1e-5), # Very low learning rate
              loss=keras.losses.BinaryCrossentropy(from_logits=True),
              metrics=[keras.metrics.BinaryAccuracy()])
```

```
# Train end-to-end. Be careful to stop before you overfit!  
model.fit(new_dataset, epochs=10, callbacks=..., validation_data=...)
```

Important note about `compile()` and `trainable`

Calling `compile()` on a model is meant to "freeze" the behavior of that model. This implies that the `trainable` attribute values at the time the model is compiled should be preserved throughout the lifetime of that model, until `compile` is called again. Hence, if you change any `trainable` value, make sure to call `compile()` again on your model for your changes to be taken into account.

Important notes about `BatchNormalization` layer

Many image models contain `BatchNormalization` layers. That layer is a special case on every imaginable count. Here are a few things to keep in mind.

- `BatchNormalization` contains 2 non-trainable weights that get updated during training. These are the variables tracking the mean and variance of the inputs.
- When you set `bn_layer.trainable = False`, the `BatchNormalization` layer will run in inference mode, and will not update its mean & variance statistics. This is not the case for other layers in general, as [weight trainability & inference/training modes are two orthogonal concepts](#). But the two are tied in the case of the `BatchNormalization` layer.
- When you unfreeze a model that contains `BatchNormalization` layers in order to do fine-tuning, you should keep the `BatchNormalization` layers in inference mode by passing `training=False` when calling the base model. Otherwise the updates applied to the non-trainable weights will suddenly destroy what the model has learned.

You'll see this pattern in action in the end-to-end example at the end of this guide.

✓ An end-to-end example: fine-tuning an image classification model on a cats vs. dogs dataset

To solidify these concepts, let's walk you through a concrete end-to-end transfer learning & fine-tuning example. We will load the Xception model, pre-trained on ImageNet, and use it on the Kaggle "cats vs. dogs" classification dataset.

✓ Getting the data

First, let's fetch the cats vs. dogs dataset using TFDS. If you have your own dataset, you'll probably want to use the utility `keras.utils.image_dataset_from_directory` to generate similar labeled dataset objects from a set of images on disk filed into class-specific

folders.

Transfer learning is most useful when working with very small datasets. To keep our dataset small, we will use 40% of the original training data (25,000 images) for training, 10% for validation, and 10% for testing.

```
1 tfds.disable_progress_bar()
2
3 train_ds, validation_ds, test_ds = tfds.load(
4     "cats_vs_dogs",
5     # Reserve 10% for validation and 10% for test
6     split=["train[:40%]", "train[40%:50%]", "train[50%:60%]"],
7     as_supervised=True, # Include labels
8 )
9
10 print(f"Number of training samples: {train_ds.cardinality()}")
11 print(f"Number of validation samples: {validation_ds.cardinality()}")
12 print(f"Number of test samples: {test_ds.cardinality()}")
```

These are the first 9 images in the training dataset -- as you can see, they're all different sizes.

```
1 plt.figure(figsize=(10, 10))
2 for i, (image, label) in enumerate(train_ds.take(9)):
3     ax = plt.subplot(3, 3, i + 1)
4     plt.imshow(image)
5     plt.title(int(label))
6     plt.axis("off")
```

We can also see that label 1 is "dog" and label 0 is "cat".

✓ Standardizing the data

Our raw images have a variety of sizes. In addition, each pixel consists of 3 integer values between 0 and 255 (RGB level values). This isn't a great fit for feeding a neural network. We need to do 2 things:

- Standardize to a fixed image size. We pick 150x150.
- Normalize pixel values between -1 and 1. We'll do this using a `Normalization` layer as part of the model itself.

In general, it's a good practice to develop models that take raw data as input, as opposed to models that take already-preprocessed data. The reason being that, if your model expects preprocessed data, any time you export your model to use it elsewhere (in a web browser, in a mobile app), you'll need to reimplement the exact same preprocessing pipeline. This gets very tricky very quickly. So we should do the least possible amount of preprocessing before hitting the model.

Here, we'll do image resizing in the data pipeline (because a deep neural network can only process contiguous batches of data), and we'll do the input value scaling as part of the model, when we create it.

Let's resize images to 150x150:

```
1 resize_fn = keras.layers.Resizing(150, 150)
2
3 train_ds = train_ds.map(lambda x, y: (resize_fn(x), y))
4 validation_ds = validation_ds.map(lambda x, y: (resize_fn(x), y))
5 test_ds = test_ds.map(lambda x, y: (resize_fn(x), y))
```

✓ Using random data augmentation

When you don't have a large image dataset, it's a good practice to artificially introduce sample diversity by applying random yet realistic transformations to the training images, such as random horizontal flipping or small random rotations. This helps expose the model to different aspects of the training data while slowing down overfitting.

```
1 augmentation_layers = [
2     layers.RandomFlip("horizontal"),
3     layers.RandomRotation(0.1),
4 ]
5
6
7 def data_augmentation(x):
8     for layer in augmentation_layers:
9         x = layer(x)
10    return x
11
12
13 train_ds = train_ds.map(lambda x, y: (data_augmentation(x), y))
```

Let's batch the data and use prefetching to optimize loading speed.

```
1 from tensorflow import data as tf_data
2
3 batch_size = 64
4
5 train_ds = train_ds.batch(batch_size).prefetch(tf_data.AUTOTUNE).cache()
6 validation_ds = validation_ds.batch(batch_size).prefetch(tf_data.AUTOTUNE).cache()
7 test_ds = test_ds.batch(batch_size).prefetch(tf_data.AUTOTUNE).cache()
```

Let's visualize what the first image of the first batch looks like after various random transformations:

```
1 for images, labels in train_ds.take(1):
2     plt.figure(figsize=(10, 10))
3     first_image = images[0]
4     for i in range(9):
5         ax = plt.subplot(3, 3, i + 1)
6         augmented_image = data_augmentation(np.expand_dims(first_image, 0))
7         plt.imshow(np.array(augmented_image[0]).astype("int32"))
8         plt.title(int(labels[0]))
9         plt.axis("off")
```

✓ Build a model

Now let's build a model that follows the blueprint we've explained earlier.

Note that:

- We add a `Rescaling` layer to scale input values (initially in the `[0, 255]` range) to the `[-1, 1]` range.
- We add a `Dropout` layer before the classification layer, for regularization.
- We make sure to pass `training=False` when calling the base model, so that it runs in inference mode, so that batchnorm statistics don't get updated even after we unfreeze the base model for fine-tuning.

```
1 base_model = keras.applications.Xception(  
2     weights="imagenet", # Load weights pre-trained on ImageNet.  
3     input_shape=(150, 150, 3),  
4     include_top=False,  
5 ) # Do not include the ImageNet classifier at the top.  
6  
7 # Freeze the base_model  
8 base_model.trainable = False  
9  
10 # Create new model on top  
11 inputs = keras.Input(shape=(150, 150, 3))  
12  
13 # Pre-trained Xception weights requires that input be scaled  
14 # from (0, 255) to a range of (-1., +1.), the rescaling layer  
15 # outputs: `(inputs * scale) + offset`  
16 scale_layer = keras.layers.Rescaling(scale=1 / 127.5, offset=-1)  
17 x = scale_layer(inputs)  
18  
19 # The base model contains batchnorm layers. We want to keep them in inference mode  
20 # when we unfreeze the base model for fine-tuning, so we make sure that the  
21 # base_model is running in inference mode here.  
22 x = base_model(x, training=False)  
23 x = keras.layers.GlobalAveragePooling2D()(x)  
24 x = keras.layers.Dropout(0.2)(x) # Regularize with dropout  
25 outputs = keras.layers.Dense(1)(x)  
26 model = keras.Model(inputs, outputs)  
27  
28 model.summary(show_trainable=True)
```

✓ Train the top layer

```
1 model.compile(  
2     optimizer=keras.optimizers.Adam(),  
3     loss=keras.losses.BinaryCrossentropy(from_logits=True),  
4     metrics=[keras.metrics.BinaryAccuracy()],  
5 )  
6  
7 epochs = 2  
8 print("Fitting the top layer of the model")  
9 model.fit(train_ds, epochs=epochs, validation_data=validation_ds)
```

✓ Do a round of fine-tuning of the entire model

Finally, let's unfreeze the base model and train the entire model end-to-end with a low learning rate.

Importantly, although the base model becomes trainable, it is still running in inference mode since we passed `training=False` when calling it when we built the model. This means that the batch normalization layers inside won't update their batch statistics. If they did, they would wreck havoc on the representations learned by the model so far.

```
1 # Unfreeze the base_model. Note that it keeps running in inference mode
2 # since we passed `training=False` when calling it. This means that
3 # the batchnorm layers will not update their batch statistics.
4 # This prevents the batchnorm layers from undoing all the training
5 # we've done so far.
6 base_model.trainable = True
7 model.summary(show_trainable=True)
8
9 model.compile(
10     optimizer=keras.optimizers.Adam(1e-5), # Low learning rate
11     loss=keras.losses.BinaryCrossentropy(from_logits=True),
12     metrics=[keras.metrics.BinaryAccuracy()],
13 )
14
15 epochs = 1
16 print("Fitting the end-to-end model")
17 model.fit(train_ds, epochs=epochs, validation_data=validation_ds)
```

After 10 epochs, fine-tuning gains us a nice improvement here. Let's evaluate the model on the test dataset:

```
1 print("Test dataset evaluation")
2 model.evaluate(test_ds)

1 for i in range(100):
2     print()
```

