

OAT40 Day 3: Transfer learning

In [1]:

```
#no_copyright_and_no_name
name = "Juan Nicolas"
surname = "Mendoza Romancio"
```

Setup

In [4]:

```
import numpy as np
import tensorflow as tf
from tensorflow import keras
```

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 raccoons may be useful to kick-start a model meant to identify tanks.

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 pre-trained 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.

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)

In [5]:

```
layer = keras.layers.Dense(3)
layer.build((None, 4)) # Create the weights
print("tensor:", len(layer.weights))
print("weight shape:", layer.weights[0].shape)
print("trainable tensors:", len(layer.trainable_weights))
print("non-trainable tensors:", len(layer.non_trainable_weights))
```

tensors: 2

weight shape: (4, 3)

trainable tensors: 2

non-trainable tensors: 0

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

In [6]:

```
layer = keras.layers.BatchNormalization()
layer.build((None, 4)) # Create the weights
print("tensor:", len(layer.weights))
print("weight shape:", layer.weights[0].shape)
print("trainable tensors:", len(layer.trainable_weights))
print("non-trainable tensors:", len(layer.non_trainable_weights))
```

tensors: 4

trainable tensors: 2

non-trainable tensors: 2

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 through training with `f1t()` or when training with any custom loop that relies on `trainable_weights` to apply gradient updates).

Example: setting `trainable` to `False`

In [7]:

```
layer = keras.layers.Dense(3)
layer.build((None, 4)) # Create the weights
layer.trainable = False

print("tensor:", len(layer.weights))
print("trainable tensors:", len(layer.trainable_weights))
print("non-trainable tensors:", len(layer.non_trainable_weights))
```

tensors: 2

trainable tensors: 0

non-trainable tensors: 2

When a trainable weight becomes non-trainable, its value is no longer updated during training.

In [8]:

```
# Make a model with 2 layers
layer1 = keras.layers.Dense(3, activation="relu")
layer2 = keras.layers.Dense(3, activation="sigmoid")
model = keras.Sequential([keras.layers.Input(shape=(3,)), layer1, layer2])
```

Freeze the first layer
layer1.trainable = False

Keep a copy of the weights of layer1 for later reference
initial_layer1_weights_values = layer1.get_weights()

Train the model
model.compile(optimizer="adam", loss="mse")

```
model.fit(np.random.random((2, 3)), np.random.random((2, 3)))
```

Check that the weights of layer1 have not changed during training
final_layer1_weights_values = layer1.get_weights()

initial_layer1_weights_values[0] == final_layer1_weights_values[0]

np.testing.assert_allclose(initial_layer1_weights_values[1], final_layer1_weights_values[1])

initial_layer1_weights_in_layer1_before_training ...

print(initial_layer1_weights_in_layer1_before_training)

initial_layer1_weights_in_layer1_after_training ...

print(initial_layer1_weights_in_layer1_after_training)

initial_layer1_weights_in_layer1_before_training ...

print(initial_layer1_weights_in_layer1_before_training)

initial_layer1_weights_in_layer1_after