

Image segmentation and neural network quantization

Here are all the import statements needed for all the exercises

```
%tensorflow_version 2.x
import tensorflow as tf
%load_ext tensorboard
import numpy as np
import matplotlib.pyplot as plt

! pip install tensorflow_model_optimization
import tensorflow_model_optimization as tfmot
```

Colab only includes TensorFlow 2.x; %tensorflow\_version has no effect. The tensorboard extension is already loaded. To reload it, use:  
%reload\_ext tensorboard  
Requirement already satisfied: tensorflow\_model\_optimization in /usr/local/lib/python3.11/dist-packages (0.8.0)  
Requirement already satisfied: absl-py~=1.2 in /usr/local/lib/python3.11/dist-packages (from tensorflow\_model\_optimization) (1.4.0)  
Requirement already satisfied: dm-tree~=0.1.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow\_model\_optimization) (0.1.8)  
Requirement already satisfied: numpy~=1.23 in /usr/local/lib/python3.11/dist-packages (from tensorflow\_model\_optimization) (1.26.4)  
Requirement already satisfied: six~=1.14 in /usr/local/lib/python3.11/dist-packages (from tensorflow\_model\_optimization) (1.17.0)

Loading the Sentinel 2 dataset. Images and labels are padded to be 256x256 in size and normalized by their maximum value. 40 images are used for the train partition (X\_train, Y\_train) and 10 for testing (X\_test, Y\_test)

```
loaded = np.load('sentinel2.npz')
X = loaded['X'].astype(np.float32) # images dataset
Y = loaded['Y'].astype(np.float32) # class dataset associated to each pixel (1 means coltivated land, 0 otherwise)
X = np.pad(X, ((0,0),(3,3),(3,3),(0,0)))
Y = np.pad(Y, ((0,0),(3,3),(3,3),(0,0)))
X_train = X[:40]/np.max(X[:40]) # dividedes datasets into training and testing
X_test = X[40:]/np.max(X[40:])
Y_train = Y[:40]
Y_test = Y[40:]
```

```
# from google.colab import drive
# drive.mount('/content/drive')
```

[TODO] Implement the U-net neural network for segmentation as drawn in the lab document.

# The structure of the U-net NN follows the one that we have seen in the lecture. In the lecture we have seen the Res-Net NN that introduces skip connections between layers.  
# This structure is useful cause the NN can better understand the relations between inputs and outputs, and it can perform channel attention or spatial attention.  
# Channel attention consists in weighting the feature map with a coefficient evaluated on the input of the layer.  
# Spatial attention is very similar and consists in weighting the pixels shared on the features of the image.  
# When we introduce skip connection we need to concatenate two output together on the features space, for doing so we must consider that they should have the same shape HxW.  
# In this NN we see also how encoding and decoding works. Encoding is used to perform features extraction through convolutional layers that reduce the size of the image but in  
# of feature maps. In this case for reducing the size of the image we are changing the strides. Strides means the jump that the kernel does on the image during the convolution  
# value equals to 2 means that we are halving the size of the input image.  
# Decondig is used to recontruct the image from the features extracted through the encoding. The decoding operation is used also for images segmentation.  
# In this case the decoding operation recontruct the image through concatenation with some encoding outputs for adding new features and so more details and then it uses UpSar  
# for increasing the size of the image. This is done by adding zeros around the image and then interpolate them with the other pixels.

```
def unet(input_shape):
    # first layer
    inputs = tf.keras.Input(shape=input_shape)
    x = tf.keras.layers.Conv2D(64, 3, strides=1, activation='relu', padding='same')(inputs) # more we go deeper on the levels and more we increase the number of feature maps ext
    x = tf.keras.layers.BatchNormalization()(x) # batch normalization performs a standardization of the input of the layer by using statistics evaluated on a batch. This is usef
    x = tf.keras.layers.ReLU()(x) # non linear activation function, important cause a linear activation function would make the NN a linear model and so we would lost the advant
    # the output of this layer is H x W x 64
    # second layer - concatenated to the first
    x = tf.keras.layers.Conv2D(64, 3, strides=1, activation='relu', padding='same')(x)
    x = tf.keras.layers.BatchNormalization()(x)
    x = tf.keras.layers.ReLU()(x)
    # the output of this layer is H x W x 64
    # third layer
    x1 = tf.keras.layers.Conv2D(128, 3, strides=2, activation='relu', padding='same')(x) # strides = 2 means that the kernel of the convolution will jump 2 pixels at a time on
    x1 = tf.keras.layers.BatchNormalization()(x1)
    x1 = tf.keras.layers.ReLU()(x1)
    # the output of this layer is H/2 x W/2 x 128
    # fourth layer
    x2 = tf.keras.layers.Conv2D(256, 3, strides=2, activation='relu', padding='same')(x1)
    x2 = tf.keras.layers.BatchNormalization()(x2)
    x2 = tf.keras.layers.ReLU()(x2)
    # the output of this layer is H/4 x W/4 x 256
    # fifth layer - concatenated to the fourth
    x2 = tf.keras.layers.Conv2D(256, 3, strides=1, activation='relu', padding='same')(x2)
    x2 = tf.keras.layers.BatchNormalization()(x2)
    x2 = tf.keras.layers.ReLU()(x2)
    x2 = tf.keras.layers.UpSampling2D(size=(2,2))(x2) # upsample the image by doubling both the dimensions
    x2 = tf.keras.layers.concatenate([x2, x1]) # concatenate the output of the layer with the output of the third layer by the features space, they should have the same image s
    # the output of this layer is H/2 x W/2 x 256+128
    # sixth layer
    x3 = tf.keras.layers.Conv2D(128, 3, strides=1, activation='relu', padding='same')(x2)
    x3 = tf.keras.layers.BatchNormalization()(x3)
    x3 = tf.keras.layers.ReLU()(x3)
    x3 = tf.keras.layers.UpSampling2D(size=(2,2))(x3)
    x3 = tf.keras.layers.concatenate([x3, x]) # concatenate the output of the layer with the output of the first-second layer by the features space, they should have the same im
    # the output of this layer is H x W x 128+64
    # seventh layer
    x4 = tf.keras.layers.Conv2D(64, 3, strides=1, activation='relu', padding='same')(x3)
    x4 = tf.keras.layers.BatchNormalization()(x4)
    x4 = tf.keras.layers.ReLU()(x4)
    # the output of this layer is H x W x 64

    # note that the last layer is not a Dense layer that performs the classification. This because we are performing image segmentation and not image classification, so the resu
    outputs = tf.keras.layers.Conv2D(2, 1, activation='softmax')(x4) # softmax will create a pdf between the two possible classes, dog or cat
    #outputs = tf.keras.layers.Conv2D(1, 1, activation='sigmoid')(x4) # sigmoid maps the input value in a probability between 0 and 1, it doesn't build a pdf. It return the clas
    model = tf.keras.Model(inputs=inputs, outputs=outputs)
    return model
```

```
Unet_model = unet((256,256,12)) # the input shape is 256x256x12
# Unet_model.summary()
```

[TODO] Compile and train the model (might take some time...)

```
Unet_model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.01),
                    loss='sparse_categorical_crossentropy',
                    metrics=['accuracy'])
Unet_model.fit(X_train, Y_train, epochs=50)
```

```
Epoch 1/50
2/2 [=====] - 34s 5s/step - loss: 1.2705 - accuracy: 0.4222
Epoch 2/50
2/2 [=====] - 1s 277ms/step - loss: 0.6016 - accuracy: 0.7847
Epoch 3/50
2/2 [=====] - 1s 278ms/step - loss: 0.5581 - accuracy: 0.8064
Epoch 4/50
2/2 [=====] - 1s 277ms/step - loss: 0.5075 - accuracy: 0.8191
Epoch 5/50
2/2 [=====] - 1s 279ms/step - loss: 0.4619 - accuracy: 0.8261
Epoch 6/50
2/2 [=====] - 1s 278ms/step - loss: 0.4242 - accuracy: 0.8262
Epoch 7/50
2/2 [=====] - 1s 278ms/step - loss: 0.4091 - accuracy: 0.8294
Epoch 8/50
2/2 [=====] - 1s 280ms/step - loss: 0.3784 - accuracy: 0.8305
Epoch 9/50
2/2 [=====] - 1s 291ms/step - loss: 0.3756 - accuracy: 0.8323
Epoch 10/50
2/2 [=====] - 1s 284ms/step - loss: 0.3703 - accuracy: 0.8348
Epoch 11/50
2/2 [=====] - 1s 277ms/step - loss: 0.3589 - accuracy: 0.8374
Epoch 12/50
2/2 [=====] - 1s 278ms/step - loss: 0.3479 - accuracy: 0.8405
Epoch 13/50
2/2 [=====] - 1s 282ms/step - loss: 0.3389 - accuracy: 0.8423
Epoch 14/50
2/2 [=====] - 1s 281ms/step - loss: 0.3345 - accuracy: 0.8421
Epoch 15/50
2/2 [=====] - 1s 278ms/step - loss: 0.3481 - accuracy: 0.8454
Epoch 16/50
2/2 [=====] - 1s 279ms/step - loss: 0.3277 - accuracy: 0.8471
Epoch 17/50
2/2 [=====] - 1s 284ms/step - loss: 0.3247 - accuracy: 0.8482
Epoch 18/50
2/2 [=====] - 1s 289ms/step - loss: 0.3193 - accuracy: 0.8490
Epoch 19/50
2/2 [=====] - 1s 291ms/step - loss: 0.3052 - accuracy: 0.8521
Epoch 20/50
2/2 [=====] - 1s 285ms/step - loss: 0.3069 - accuracy: 0.8521
Epoch 21/50
2/2 [=====] - 1s 285ms/step - loss: 0.3002 - accuracy: 0.8569
Epoch 22/50
2/2 [=====] - 1s 285ms/step - loss: 0.3021 - accuracy: 0.8572
Epoch 23/50
2/2 [=====] - 1s 286ms/step - loss: 0.2965 - accuracy: 0.8619
Epoch 24/50
2/2 [=====] - 2s 291ms/step - loss: 0.2858 - accuracy: 0.8675
Epoch 25/50
2/2 [=====] - 1s 288ms/step - loss: 0.2851 - accuracy: 0.8707
Epoch 26/50
2/2 [=====] - 1s 281ms/step - loss: 0.2759 - accuracy: 0.8729
Epoch 27/50
2/2 [=====] - 1s 296ms/step - loss: 0.2761 - accuracy: 0.8767
Epoch 28/50
2/2 [=====] - 1s 298ms/step - loss: 0.2775 - accuracy: 0.8740
Epoch 29/50
2/2 [=====] - 1s 300ms/step - loss: 0.2790 - accuracy: 0.8767
```

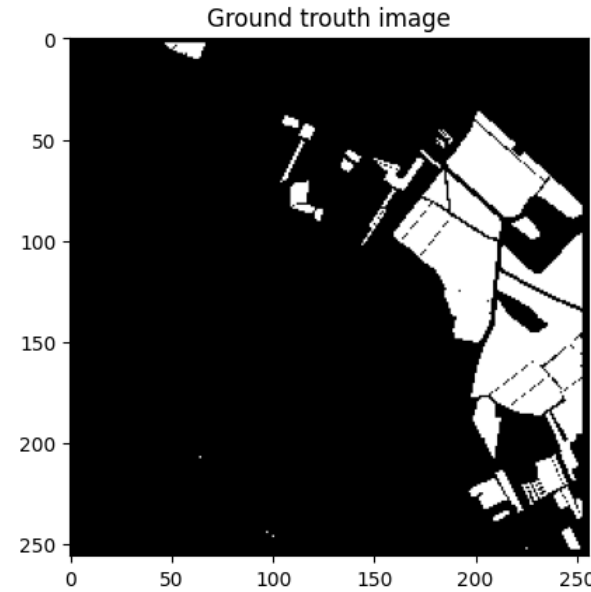
[TODO] Test the model on the test set and measure the accuracy.

```
Unet_model.evaluate(X_test, Y_test)

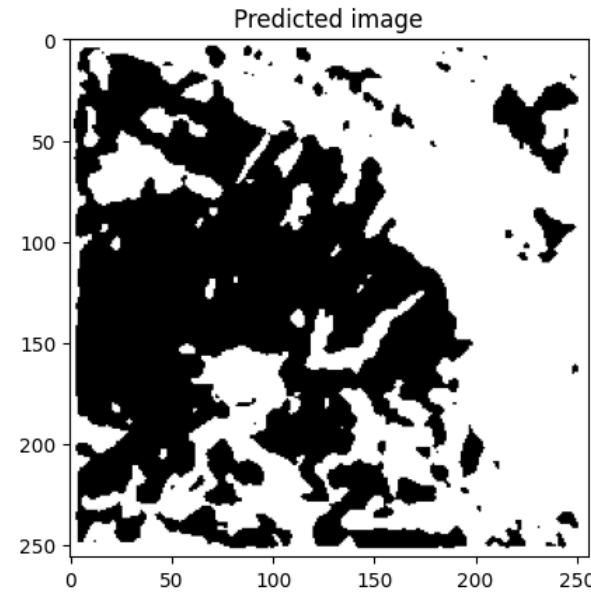
image = Y_test[1]
plt.imshow(image, 'gray')
plt.title('Ground trouth image')
plt.show()

predict_im = Unet_model.predict(X_test[1].reshape(1,256,256,12))
predict_im = np.squeeze(predict_im)
predict_im = np.argmax(predict_im, axis=-1)
plt.imshow(predict_im, 'gray')
plt.title('Predicted image')
plt.show()
```

1/1 [=====] - 2s 2s/step - loss: 0.9591 - accuracy: 0.5645



1/1 [=====] - 1s 599ms/step



[TO DO] Convert model to TFLite with 8-bit weight quantization

# in our model we have many parameters and each is saved as float32. This cost to use a lot of memory space, so we want to reduce the dimension of these parameters. For doing  
# perform quantization, so we can quantize each float parameter to an int parameter that can be rappresented with only 8 bits.  
# PTQ (Post Training Quantization) quantize each parameter on an already trained neural network. This type of quantization is the faster and the easiest to implement, but it  
# accuracy of the model, cause the quantization is chaning the optimal parameters that the model has found during the training.

```
converter = tf.lite.TFLiteConverter.from_keras_model(Unet_model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]
```

```
quantized_model = converter.convert()
with open("quantized_model.tflite", "wb") as f:
    f.write(quantized_model)
```

[TODO] Test the accuracy of the quantized model by writing your own "evaluate" function. Remember that TFLite interpreter can only process one sample at a time, not a batch.

```
# Initialize the interpreter
interpreter = tf.lite.Interpreter('quantized_model.tflite')
interpreter.allocate_tensors()

input_details = interpreter.get_input_details()
output_details = interpreter.get_output_details()

def evaluate(interpreter, X_test, Y_test):
    accuracies = []

    for i, (input_data, true_class) in enumerate(zip(X_test, Y_test)):
        input_data = input_data.astype(input_details[0]['dtype']) # read the image from the test dataset

        interpreter.set_tensor(input_details[0]['index'], np.expand_dims(input_data, axis=0)) # set the input tensor with the image and add a dimension for the batch

        interpreter.invoke() # run the inference on the image
        output_data = interpreter.get_tensor(output_details[0]['index']) # get the output tensor of the model
        predicted_class = np.squeeze(output_data) # remove the batch dimension
        predicted_class = np.argmax(predicted_class, axis=2) # get the class with the highest probability

        true_class = np.squeeze(true_class) # remove the batch dimension
        accuracy = np.mean(predicted_class == true_class) # evaluate the accuracy of the model
        accuracies.append(accuracy) # save the accuracy of the model

    print(f"Image {i+1} - Accuracy: {accuracy}") # print the accuracy of the model

    plt.imshow(true_class, 'gray') # show the image
    plt.title('Ground trouth image')
    plt.show()

    plt.imshow(predicted_class, 'gray') # show the predicted image
    plt.title('Predicted image')
    plt.show()
    return np.mean(accuracies)
# Call the evaluate function
accuracy = evaluate(interpreter, X_test, Y_test)
print(f"Mean accuracy: {accuracy}")

original_predictions = Unet_model.predict(X_test)
original_predictions = (original_predictions > 0.5).astype(np.uint8)
```

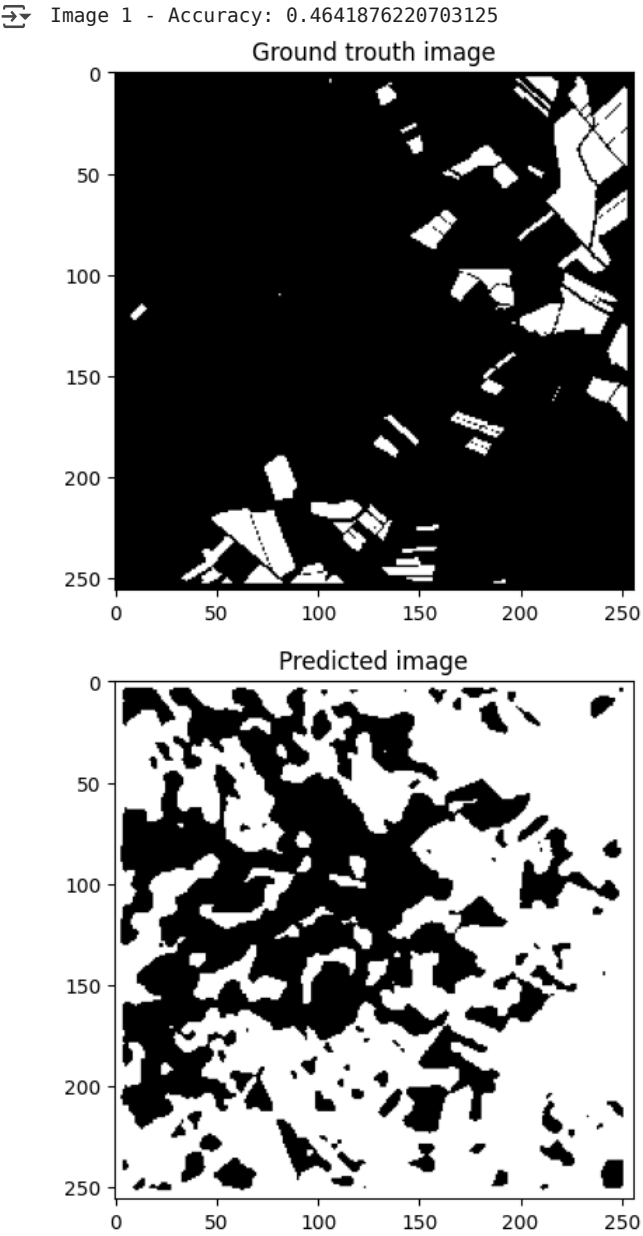


Image 2 - Accuracy: 0.5809478759765625

Ground trouth image

**[TODO]** Finetune the Keras model using quantization-aware training and measure the accuracy on the test set after actually quantizing it

```
# the best way to quantize the parameter is through QAT (Quantization Aware Training). This method consists in training the model to find the best parameters that can be quantized.
# So during the training the model will be aware that the parameters need to be quantized, so it will find the best values to be quantized.
# A model quantized with QAT can perform similar to the original model, but with a lower memory occupation. Note that QAT model need an ad-hoc training and calibration.

# cannot apply QAT to batch normalizer layers
quant_aware_model = tfmot.quantization.keras.quantize_model(Unet_model)
quant_aware_model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.01),
                          loss='sparse_categorical_crossentropy',
                          metrics=['accuracy'])
quant_aware_model.fit(X_train, Y_train, epochs=50)

converter = tf.lite.TFLiteConverter.from_keras_model(quant_aware_model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]
quantized_aware_model = converter.convert()
```

↗

```
-----
RuntimeError                                Traceback (most recent call last)
<ipython-input-12-7f43f456f466> in <cell line: 6>()
      4
      5 # cannot apply QAT to batch normalizer layers
----> 6 quant_aware_model = tfmot.quantization.keras.quantize_model(Unet_model)
      7 quant_aware_model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.01),
      8                           loss='sparse_categorical_crossentropy',
```