## Breast cancer classification with Keras and Deep Learning **Import Packages** In $[1]\colon$ $\mid$ # set the matplotlib backend so figures can be saved in the background import matplotlib matplotlib.use("Agg") # import the necessary packages from tensorflow.keras.preprocessing.image import ImageDataGenerator from tensorflow.keras.callbacks import LearningRateScheduler from tensorflow.keras.models import Sequential from tensorflow.keras.layers import BatchNormalization from tensorflow.keras.layers import SeparableConv2D from tensorflow.keras.layers import MaxPooling2D from tensorflow.keras.layers import Activation from tensorflow.keras.layers import Flatten from tensorflow.keras.layers import Dropout from tensorflow.keras.layers import Dense from tensorflow.keras.optimizers import Adagrad from tensorflow.keras import backend as K from tensorflow.keras.utils import to\_categorical from imutils import paths import matplotlib.pyplot as plt import numpy as np import argparse import random import shutil import os Our Config class In [2]: | class Config: # initialize the path to the \*original\* input directory of images ORIG\_INPUT\_DATASET = "breast-cancer-classification/data" # initialize the base path to the \*new\* directory that will contain # our images after computing the training and testing split BASE\_PATH = "breast-cancer-classification/data" # derive the training, validation, and testing directories TRAIN\_PATH = os.path.sep.join([BASE\_PATH, "training"]) VAL\_PATH = os.path.sep.join([BASE\_PATH, "validation"]) TEST\_PATH = os.path.sep.join([BASE\_PATH, "testing"]) # define the amount of data that will be used training TRAIN\_SPLIT = 0.8 # the amount of validation data will be a percentage of the # \*training\* data VAL\_SPLIT = 0.1 # initialize our config class config = Config() Building the breast cancer image dataset In [3]: # grab the paths to all input images in the original input directory # and shuffle them imagePaths = list(paths.list\_images(config.ORIG\_INPUT\_DATASET)) random.seed(42) random.shuffle(imagePaths) # compute the training and testing split i = int(len(imagePaths) \* config.TRAIN\_SPLIT) trainPaths = imagePaths[:i] testPaths = imagePaths[i:] # we'll be using part of the training data for validation i = int(len(trainPaths) \* config.VAL\_SPLIT) valPaths = trainPaths[:i] trainPaths = trainPaths[i:] # define the datasets that we'll be building ("training", trainPaths, config.TRAIN\_PATH), ("validation", valPaths, config.VAL\_PATH), ("testing", testPaths, config.TEST\_PATH) In [4]: # Loop over the datasets for (dType, imagePaths, baseOutput) in datasets: # show which data split we are creating print("[INFO] building '{}' split".format(dType)) # if the output base output directory does not exist, create it if not os.path.exists(baseOutput): print("[INFO] 'creating {}' directory".format(baseOutput)) os.makedirs(baseOutput) # loop over the input image paths for inputPath in imagePaths: # extract the filename of the input image and extract the # class label ("0" for "negative" and "1" for "positive") filename = inputPath.split(os.path.sep)[-1] label = filename[-5:-4]# build the path to the label directory labelPath = os.path.sep.join([baseOutput, label]) # if the label output directory does not exist, create it if not os.path.exists(labelPath): print("[INFO] 'creating {}' directory".format(labelPath)) os.makedirs(labelPath) # construct the path to the destination image and then copy # the image itself p = os.path.sep.join([labelPath, filename]) shutil.copy2(inputPath, p) [INFO] building 'training' split \_\_\_\_\_ SameFileError Traceback (most recent call last) <ipython-input-4-e794fee51c6c> in <module> # the image itself p = os.path.sep.join([labelPath, filename]) ---> 30 shutil.copy2(inputPath, p) c:\users\leobr\appdata\local\programs\python\python39\lib\shutil.py in copy2(src, dst, follow\_symlinks) 433 **if** os.path.isdir(dst): dst = os.path.join(dst, os.path.basename(src)) --> 435 copyfile(src, dst, follow\_symlinks=follow\_symlinks) copystat(src, dst, follow\_symlinks=follow\_symlinks) c:\users\leobr\appdata\local\programs\python\python39\lib\shutil.py in copyfile(src, dst, follow\_symlinks) if \_samefile(src, dst): 243 raise SameFileError("{!r} and {!r} are the same file".format(src, dst)) --> 244 246 file\_size = 0 SameFileError: 'breast-cancer-classification/data\\training\\0\\9036\_idx5\_x551\_y2201\_class0.png' and 'breast-cancer-classification/data\\training\\0\\9036\_idx5\_x551\_y2201\_class0.png' are the same file CancerNet: Our breast cancer prediction CNN In [5]: class CancerNet: @staticmethod def build(width, height, depth, classes): # initialize the model along with the input shape to be # "channels last" and the channels dimension itself model = Sequential() inputShape = (height, width, depth) chanDim = -1# if we are using "channels first", update the input shape # and channels dimension if K.image\_data\_format() == "channels\_first": inputShape = (depth, height, width) chanDim = 1# CONV => RELU => POOL model.add(SeparableConv2D(32, (3, 3), padding="same", input\_shape=inputShape)) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(MaxPooling2D(pool\_size=(2, 2))) model.add(Dropout(0.25)) # (CONV => RELU => POOL) \* 2 model.add(SeparableConv2D(64, (3, 3), padding="same")) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(SeparableConv2D(64, (3, 3), padding="same")) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(MaxPooling2D(pool\_size=(2, 2))) model.add(Dropout(0.25)) # (CONV => RELU => POOL) \* 3 model.add(SeparableConv2D(128, (3, 3), padding="same")) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(SeparableConv2D(128, (3, 3), padding="same")) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(SeparableConv2D(128, (3, 3), padding="same")) model.add(Activation("relu")) model.add(BatchNormalization(axis=chanDim)) model.add(MaxPooling2D(pool\_size=(2, 2))) model.add(Dropout(0.25)) # first (and only) set of FC => RELU layers model.add(Flatten()) model.add(Dense(256)) model.add(Activation("relu")) model.add(BatchNormalization()) model.add(Dropout(0.5)) # softmax classifier model.add(Dense(classes)) model.add(Activation("softmax")) # return the constructed network architecture return model Our training script In [6]: # # construct the argument parser and parse the arguments # ap = argparse.ArgumentParser() # ap.add\_argument("-p", "--plot", type=str, default="plot.png", # help="path to output loss/accuracy plot") # args = vars(ap.parse\_args()) # since we are using Jupyter Notebooks we can replace our argument # parsing code with \*hard coded\* arguments and values args **=** { "plot": "plot.png", In [7]: # initialize our number of epochs, initial learning rate, and batch $NUM_EPOCHS = 40$ $INIT_LR = 1e-2$ BS = 32# determine the total number of image paths in training, validation, # and testing directories trainPaths = list(paths.list\_images(config.TRAIN\_PATH)) totalTrain = len(trainPaths) totalVal = len(list(paths.list\_images(config.VAL\_PATH))) totalTest = len(list(paths.list\_images(config.TEST\_PATH))) # calculate the total number of training images in each class and # initialize a dictionary to store the class weights trainLabels = [int(p.split(os.path.sep)[-2]) for p in trainPaths] trainLabels = to\_categorical(trainLabels) classTotals = trainLabels.sum(axis=0) classWeight = dict() # loop over all classes and calculate the class weight for i in range(0, len(classTotals)): classWeight[i] = classTotals.max() / classTotals[i] In [8]: # initialize the training data augmentation object trainAug = ImageDataGenerator( rescale=1 / 255.0, rotation\_range=20, zoom\_range=0.05, width\_shift\_range=0.1, height\_shift\_range=0.1, shear\_range=0.05, horizontal\_flip=True, vertical\_flip=True, fill\_mode="nearest") # initialize the validation (and testing) data augmentation object valAug = ImageDataGenerator(rescale=1 / 255.0)

In [9]: # initialize the training generator

trainGen = trainAug.flow\_from\_directory(

class\_mode="categorical",
target\_size=(48, 48),
color\_mode="rgb",
shuffle=True,
batch\_size=BS)

# initialize the validation generator

config.TRAIN\_PATH,

```
valGen = valAug.flow_from_directory(
       config.VAL_PATH,
       class_mode="categorical",
       target_size=(48, 48),
       color_mode="rgb",
       shuffle=False,
       batch_size=BS)
    # initialize the testing generator
    testGen = valAug.flow_from_directory(
       config.TEST_PATH,
       class_mode="categorical",
       target_size=(48, 48),
       color_mode="rgb",
       shuffle=False,
       batch_size=BS)
    Found 255556 images belonging to 2 classes.
   Found 42596 images belonging to 2 classes.
   Found 99743 images belonging to 2 classes.
In [10]: | ### initialize our CancerNet model and compile it
    model = CancerNet.build(width=48, height=48, depth=3,
       classes=2)
    opt = Adagrad(lr=INIT_LR, decay=INIT_LR / NUM_EPOCHS)
    model.compile(loss="binary_crossentropy", optimizer=opt,
       metrics=["accuracy"])
    # fit the model
    H = model.fit(
       x=trainGen,
       steps_per_epoch=totalTrain // BS,
       validation_data=valGen,
       validation_steps=totalVal // BS,
       class_weight=classWeight,
       epochs=NUM_EPOCHS)
   c:\users\leobr\appdata\local\programs\python\python39\lib\site-packages\tensorflow\python\keras\optimizer_v2\optimizer_v2.py:374: UserWarning: The `lr` argument is deprecated, use `learning_rate` instead.
    WARNING:tensorflow:From c:\users\leobr\appdata\local\programs\python\python\python39\lib\site-packages\tensorflow.python.ops.array_ops) with validate_indices is deprecated and will be removed in a future version.
    The `validate_indices` argument has no effect. Indices are always validated on CPU and never validated on GPU.
   Epoch 1/40
    Epoch 2/40
    Epoch 3/40
    Epoch 4/40
    Epoch 5/40
    Epoch 6/40
    Epoch 7/40
    Epoch 8/40
    Epoch 9/40
    Epoch 10/40
    Epoch 11/40
    Epoch 12/40
    Epoch 13/40
    Epoch 14/40
    Epoch 15/40
    Epoch 16/40
    Epoch 17/40
    Epoch 18/40
    Epoch 19/40
    Epoch 20/40
    Epoch 21/40
    Epoch 22/40
    Epoch 23/40
    Epoch 24/40
    Epoch 25/40
    Epoch 26/40
    Epoch 27/40
    Epoch 28/40
    Epoch 29/40
    Epoch 30/40
    Epoch 31/40
    Epoch 32/40
    Epoch 33/40
    Epoch 34/40
    Epoch 35/40
    Epoch 36/40
    Epoch 37/40
    Epoch 38/40
    Epoch 39/40
    Epoch 40/40
    In [17]:
    import sklearn
    import sklearn.metrics
In [18]: # reset the testing generator and then use our trained model to
    # make predictions on the data
    print("[INFO] evaluating network...")
    testGen.reset()
    predIdxs = model.predict(x=testGen, steps=(totalTest // BS) + 1)
    # for each image in the testing set we need to find the index of the
    # label with corresponding largest predicted probability
    predIdxs = np.argmax(predIdxs, axis=1)
    # show a nicely formatted classification report
    print(sklearn.metrics.classification_report(testGen.classes, predIdxs,
       target_names=testGen.class_indices.keys()))
    [INFO] evaluating network...
         precision recall f1-score support
           0.93
               0.85
                   0.89
           0.69
                   0.76
                       28448
                0.84
                    0.85
     accuracy
     macro avg
           0.81
               0.84
                   0.82
                       99743
           0.86
               0.85
                   0.85
                       99743
    weighted avg
In [20]: # compute the confusion matrix and and use it to derive the raw
    # accuracy, sensitivity, and specificity
    cm = sklearn.metrics.confusion_matrix(testGen.classes, predIdxs)
    total = sum(sum(cm))
    acc = (cm[0, 0] + cm[1, 1]) / total
    sensitivity = cm[0, 0] / (cm[0, 0] + cm[0, 1])
    specificity = cm[1, 1] / (cm[1, 0] + cm[1, 1])
    # show the confusion matrix, accuracy, sensitivity, and specificity
    print(cm)
    print("acc: {:.4f}".format(acc))
    print("sensitivity: {:.4f}".format(sensitivity))
    print("specificity: {:.4f}".format(specificity))
    [[60707 10588]
    [ 4632 23816]]
    acc: 0.8474
    sensitivity: 0.8515
   specificity: 0.8372
In [27]: # plot the training loss and accuracy
    N = NUM_EPOCHS
    plt.style.use("ggplot")
    plt.figure()
    plt.plot(np.arange(0, N), H.history["loss"], label="train_loss")
    plt.plot(np.arange(0, N), H.history["val_loss"], label="val_loss")
    plt.plot(np.arange(0, N), H.history["accuracy"], label="train_acc")
    plt.plot(np.arange(0, N), H.history["val_accuracy"], label="val_acc")
    plt.title("Training Loss and Accuracy on Dataset")
    plt.xlabel("Epoch #")
    plt.ylabel("Loss/Accuracy")
    plt.legend(loc="lower left")
    plt.savefig(args["plot"])
    plt.show()
    <ipython-input-27-d02a2d202664>:14: UserWarning: Matplotlib is currently using agg, which is a non-GUI backend, so cannot show the figure.
    plt.show()
```

Code License Agreement

Copyright (c) 2020 PyImageSearch.com

LONGER, FORMAL VERSION

SIMPLE VERSION

Feel free to use this code for your own projects, whether they are purely educational, for fun, or for profit. THE EXCEPTION BEING if you are developing a course, book, or other educational product. Under \*NO CIRCUMSTANCE\* may you use this code for your own paid educational or self-promotional ventures without written consent from Adrian Rosebrock and PyImageSearch.com.

Permission is hereby granted, free of charge, to any person obtaining

a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions: The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software. Notwithstanding the foregoing, you may not use, copy, modify, merge, publish, distribute, sublicense, create a derivative work, and/or sell copies of the Software in any work that is designed, intended, or marketed for pedagogical or instructional purposes related to programming, coding, application development, or information technology. Permission for such use, copying, modification, and merger, publication, distribution, sub-licensing, creation of derivative works, or sale is expressly withheld. THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.