

OVERVIEW

PLANT CLASSIFICATION

1

• DATASET

- προβλήματα με προς την κατανομή των εικόνων. The long-tail problem.
- με βάση διαχωρισμός train / val / test. Υπάρχει ένας φακέλος με τα train data (PlantCLEF2017Train1EOL) του EOL και ένας φακέλος με τα test data (PlantCLEF2017Test), ο οποίος έχει λίγα εικόνες και όχι σε φακέλους και παρόμοια οι εικόνες αυτές δεν ήταν από διαφορετικές κλάσεις. Υπάρχει δηλαδή πρόβλημα, καθώς το train set μπορεί να έχει X κλάσεις, αλλά το test set έχει λίγες.

★ ΛΥΣΗ

Χρησιμοποιούμε τον ImageDataGenerator του TF2 για να διαχωρίσουμε τις εικόνες και αυτό για να μπορεί να γίνει σε μία ποσότητα διαδικασίας. Διατηρούμε το PlantCLEF2017Train1EOL με τις 10000 κλάσεις του. Δημιουργούμε έναν δικό μας φακέλο PlantCLEF2017TestEOL, ο οποίος έχει test δεδομένα για ΟΛΕΣ τις κλάσεις. Προκαθόρισε να γίνει αυτό, ακολουθώντας την εξής διαδικασία: Πρώτα ανόητα χρησιμοποιούμε το PlantCLEF2017OnlyTest.csv αρχείο για να βρούμε σε ένα φακέλο όλες τις εικόνες που αντιστοιχούν σε μία κλάση και παρόμοια διατηρούμε το εποχιακό της ονομασίας, ονομάζουμε αυτό το φακέλο με το όνομα της ιδίας της κλάσης. Κα... το ίδιο για όλες τις κλάσεις. Παρατηρούμε ότι έχουμε test δεδομένα για πολλές λίγες κλάσεις (π.χ. 1/5 από τις 10000 περίπου). Για να το αντιμετωπίσουμε αυτό, βρίσκουμε σε ποιές κλάσεις έχουμε test δεδομένα και προσθέτουμε από τις αντιστοιχικές κλάσεις το 10% των train δεδομένων τους με κατάδει των 1 φωτογραφία. Δηλαδή είναι περίπου που κάποια κλάση από train δεδομένα έχει λίγες από 10 φωτογραφίες (και αυτό θα είναι σε 1 στο test set θα περάσει καλά) έτσι ορίζουμε αυτή να είναι να πάρει έτσι 1 φωτογραφία για test δεδομένα της. Ο κώδικας για αυτές τις διαδικασίες βρίσκεται στο αρχείο "create_test_datASET.py", το οποίο δημιουργεί τον φάκελο με τα ΕΛΛΗΝΗ test δεδομένα χρησιμοποιώντας το αρχείο .csv και "create_all_classes_test_datASET.py", το οποίο βρίσκει σε ποιές κλάσεις έχουμε test δεδομένα και προσθέτει όπως είναι παραπάνω.

Το αρχείο `"create_generators.py"` δημιουργεί τους Image Data Generators που χρησιμοποιεί το validation set προκύπτει αυτόματα με split στο 10 train set (100% - 20% split οι καλύτερες επιλογές). Τελικά δοκιμάστηκε και η 50% ποσοστώσεως split, με το 20% να μικραίνει, ~~και να δοκιμάζεται~~

Όσον αφορά το πρόβλημα του unbalanced dataset (long-tail problem), η λύση που προτιμάμε είναι να δημιουργήσουμε fake data. Έχουν δοκιμάσει διάφορες τεχνικές (shift, rotate, flip, crop etc.) και επιλέχθηκαν αυτοί που φαίνονται στο αρχείο `"data_augmentation.py"` καθώς και οι αντίστοιχες τιμές που φαίνονται. ΔΕΝ επιλέχθηκε να εφαρμοστεί ένας keras-faceli / γιλιτρίσις τυχαία σε μία εικόνα κατά τη διαδικασία της εκπαίδευσης (on-the-fly), διατηρώντας έτσι το μέγεθος των εικόνων ίδιο, ΑΛΛΑ επιλέχθηκε να ορίσουμε έναν αριθμό (1000 εικόνες) και να εφαρμοστεί a priori τους μετασχηματισμούς ΤΟΣΟ ΠΟΣΟΥ η κάθε κλάση να έχει τον αριθμό των μετασχηματισμένων εικόνων που ορίσαμε. Αυτό έχει γίνει προγραμματίζοντας ΤΕΡΑΣΤΙΑ διαφορές στο accuracy. Με λίγα λόγια, με αυτό τον τρόπο μπορεί να οριστεί ένας αριθμός εικόνων και να εφαρμοστεί τυχαίους μετασχηματισμούς μέχρι να φτάσουμε αυτόν τον αριθμό per class. Κάτι είναι ο αριθμός αυτός να γιν είναι μεγάλος αριθμός (20 1000 είναι καλό) γιατί μετά παρατηρείται το φαινόμενο του να έχουμε των ίδιων εικόνων 500 φορές. Επειδή οι 10000 κλάσεις του PlantCLEF EoL είναι πολλές και δεν είναι στην ζώνη του φαινομένου του δικού μας προβλήματος, επιλέχθηκε 100 κλάσεις τυχαία και κινούμε augment αυτές. Δημιουργείται έτσι του φακέλου 100classesPlantCLEF2017trainEoL ο οποίος έχει μέσα του φακέλο data με τα δεδομένα χωρίς το augment και τον φακέλο augmented data που έχει τα augmented data που χρησιμοποιούμε.

Αντίστοιχα δημιουργείται ο φακέλος 100classesPlantCLEF2017testEoL που περιλαμβάνει τα test data και τις αντίστοιχες 100 κλάσεις.

ΔΙΚΤΥΑ

- γενικά η πρώτη πρόταση ήταν να δοκιμάσει τα ίδια δίκτυα που δοκιμάσε στο paper του ο vikarij του διαγωνισμού. Το πρόβλημα είναι ότι όταν δοκιμάστηκε αυτό το δίκτυο, δοκιμάστηκε με στο dataset των 10000 κτηνών και αυτό έδειξε ότι ήταν χρονοβόρο, δηλαδή περίμενε κιόλας να γίνουν οι κατάλληλες δοκιμές και να αναπαράγει fine-tune στις hyperparameters. Έτσι, παρατηρήθηκε το φαινόμενο το ~~validation~~ set accuracy να είναι στο $\approx 40\%$, ενώ το ~~training~~ accuracy έφτανε το $\approx 100\%$. Υπάρχει, δηλαδή, πρόβλημα OVERFIT. Το δίκτυο δοκιμάστηκε να εκπαιδευτούν έχοντας ως αρχικοποίηση τα βάρη του ImageNet, και να ΑΚΟΜΑ βγούνε στα τυφλά experiments. Έγινε προσέγγιση να διαπιστωθεί αν είναι καλύτερο τα βάρη να είναι frozen, semi-frozen ή non frozen. Σε γενικές γραμμές τα frozen βάρη δεν ανέδωσαν και αυτό γιατί τα high level features του CNN για το πρόβλημα της ταξινόμησης των ζώων δε μπορεί να είναι ίδια με τα extracted/learned high level features του προβλεπτή της ταξινόμησης των 10000 ΓΕΝΙΚΩΝ κτηνών του dataset. Όταν συγκρίνατε σε fine semi-frozen ή non, όπως τα low level layers είναι frozen (αυτές που εφόσον την πληροφορία των αυτών, γενικών κτηνών), ενώ τα high level layers είναι unfrozen, για να μπορεί να προσαρμοστεί στο πρόβλημα της ταξινόμησης των ζώων. Το δίκτυο αφού κάνει την εξαγωγή των χαρακτηριστικών κατέληξε σε ένα Global Average Pooling, έπειτα από το οποίο επιλέχθηκε να τονομετρηθεί για ταξινόμηση ένα Dense/Fully Connected Layer με 2048 νευρώνες και 0.5 Dropout (για να εξαλειφθεί η overfitting) και τέλος ένα τελευταίο απλό Dense/Fully Connected οι νευρώνες θα λάβουν κατά) και τέλος ένα τελευταίο απλό Dense/Fully Connected με νευρώνες 10000 δόσε και οι κλήσεις της. Η αρχικοποίηση αυτή γίνεται κοιτώντας/υποδεικνύοντας και στις λεπτομέρειες δοκιμών. Το semi-frozen βάρη ανέδωσαν καλύτερα

από ότι πριν να freeze. Τέλος, δοκιμάστηκε να γίνει unfreeze 4
2η βάρη όλα και να είναι learnable και αυτό απέδωσε καλύτερα από
0,71 αθρ.

- Έχοντας περνούσει για αρχιτεκτονική στο final, είχε παραμείνει το πρόβλημα
του αριθμού των κλάσεων και του unbalanced dataset, το οποίο συνεπαγόταν
να χάνεται πολύς χρόνος στις δοκιμές. Αποφασίστηκε να μειωθούν οι κλάσεις
και λήφουν με τον εφικτό τρόπο να κρατηθούν μόνο οι κλάσεις που είχαν πάνω από
30 samples $\rightarrow \approx 2000$ κλάσεις. Έτσι και μειώθηκε ο αριθμός των κλάσεων και έγινε κάπως
πιο balanced το dataset. Αυτό, δεν είχε κάνει augmentation και δεν είχε γίνει
το test dataset (που αναφέρεται στις βελ. 1,2). Το πρόβλημα του OVERFIT συνέχισε
να υπάρχει. Το πρόβλημα αυτό μπορεί να οφείλεται / οφείλουν στο γεγονός ότι
τα δεδομένα για κάποιες κλάσεις (που περιγράφηκε) ήταν λίγα και τα δίκτυα μάλλον
βασίζονταν να τα memorize. Διότι το δίκτυο έπρεπε ότι μάλλον να μάθει όλες
ήταν κάποιες εικόνες (εφόσον δεν είχε πολλές) και έτσι είχαν γενικευτική ικανότητα. Δοκιμάστηκε
να κρατηθούν οι κλάσεις που έχουν πάνω από 30 samples $\rightarrow \approx 1000$ κλάσεις, αλλά το
ίδιο πρόβλημα παρέμεινε. Έτσι λοιπόν, αποφασίστηκε να γίνει augmentation έτσι
ώστε να έχουμε σχετικά αριθμό εικόνων για κάθε κλάση (που περιγράφεται στις
βελτίες 1,2). Μεταβιβάστηκαν λοιπόν από το Fig. 1 στο Fig. 2:

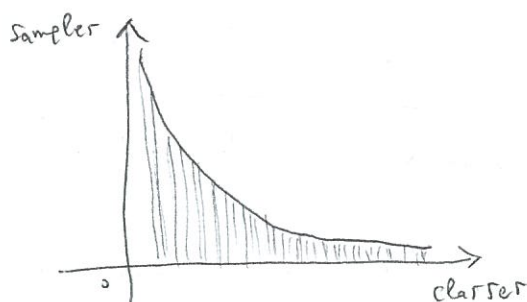


Fig. 1

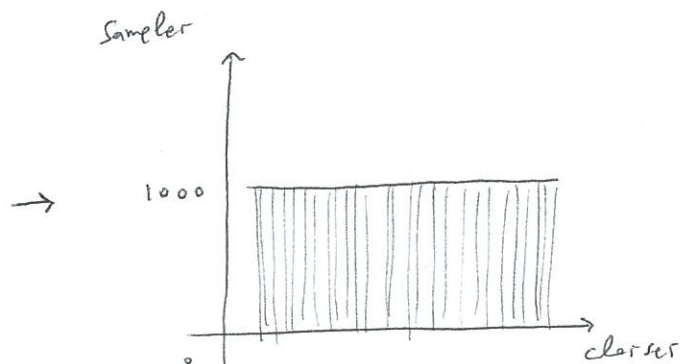
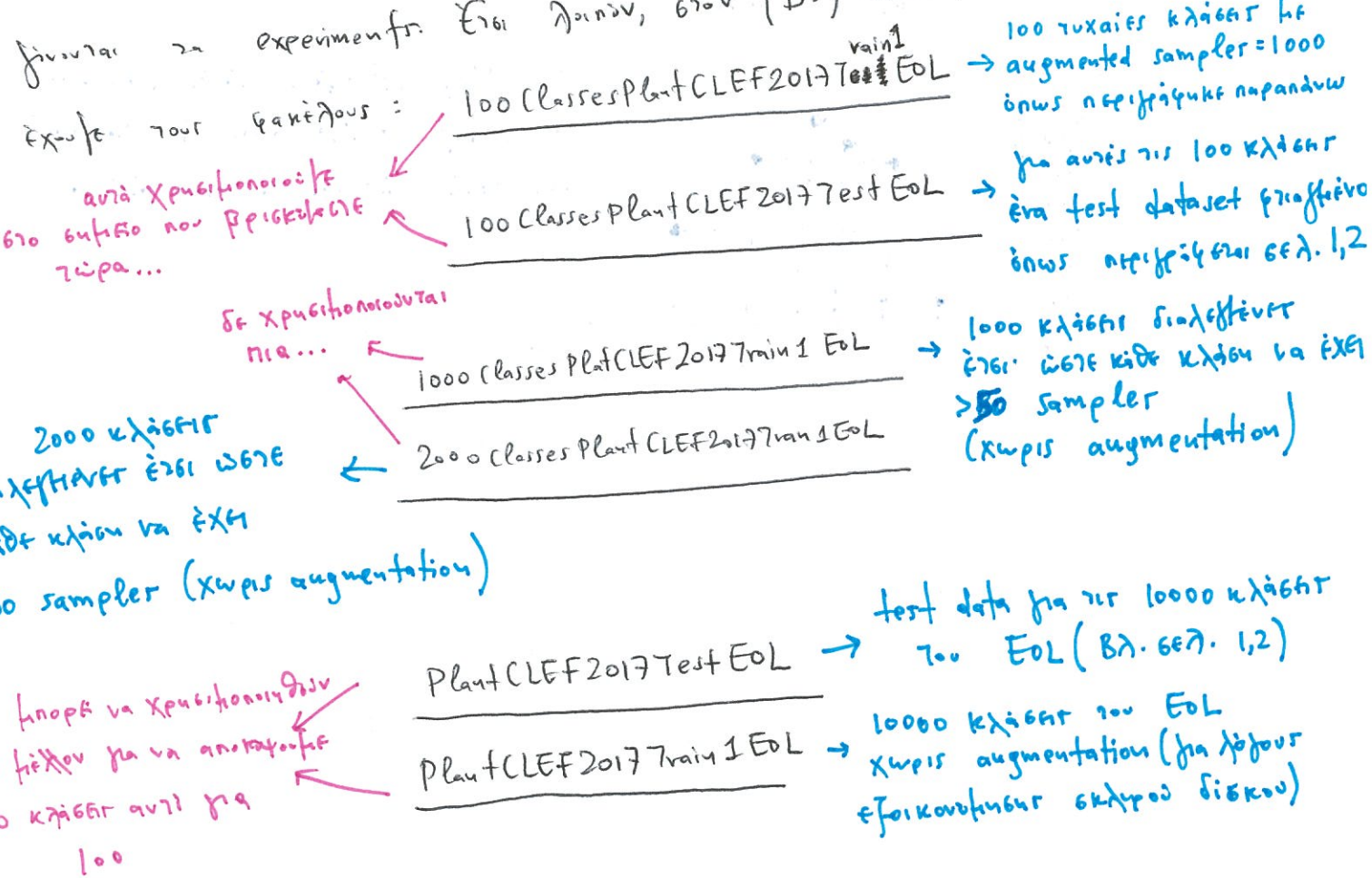


Fig. 2

Διευκρινίζεται πως (όπως φαίνεται στα Fig. 1 και Fig. 2) η augmentation 5
εγκύβει συστηματικά τυχαιούς μετασχηματισμούς, ώστε να δημιουργηθούν νέα fake
data (και έτσι να αυξηθεί ο αριθμός των sampler για κάθε κλάση και για όλα τα epochs)
Όπως αναφέρεται και στα σελ. 1,2 η augmentation δε γίνεται on-the-fly
(λόγω memory error) αλλά αποθηκεύονται οι εικόνες σε φακέλο αποθηκευμάτων η
βασική της αναφοράς. Augmentation έχει ως κύριο χαρακτηριστικό ότι
100 τυχαιές κλάσεις, καθώς η 100 σε νοίκι είναι στην βάση δεδομένων του dataset
που θα έχουμε από Map210 και είναι και ένα καλό νοίκι υπολογιστικά για να
γίνονται τα experiments. Έτσι λοιπόν, στον (D:) του laptop του πατέρα

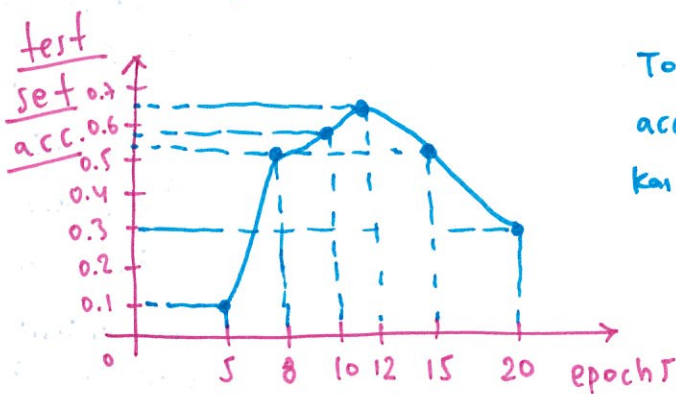


να συζητηθεί εδώ πως αν θέλουμε να αποθηκεύουμε 500 κλάσεις αυτή για 100 και να κάνουμε augment,
τότε πρώτα αν' εδω πρέπει να περάσουμε το σωστό δίκτυο και να επιβεβαιώσουμε το μέγεθος
θα λείβουν σε αυτόν και δεύτερον να χρησιμοποιήσουμε το αρχείο "data-augmentation.py",
το οποίο παίρνει ως input έναν φάκελο με X ^{κλάσεις} ~~data~~ και βγαίνει X ^{κλάσεις με} augmented data.
Μια συζήτηση ακόμα που πρέπει να γίνει και αφορά τη διαίρεση είναι ότι δημιουργώντας
fake data δε λύνουμε απλώς το πρόβλημα. Ουσιαστικά δε δημιουργούμε νέα δεδομένα,

αλλά εστιάσαμε πιο συγκεκριμένα στα υβρίδα υμείκτους και έτσι αυξήσαμε τον αριθμό τους. Και που χρειάζεται ft αυξ, είναι το γεγονός ότι χρησιμοποιώντας τα augmented data το training και validation accuracy γίνονται από ένα σύνολο το $\approx 100\%$, αλλά το test accuracy είναι αρκετά χαμηλά. Για να γίνει περισσότερο καθαρό αυτό, θα παραδείξω τις test set accuracy του NASNetMobile δικτύου, στο οποίο έχουν πολλή δομή, και αυτό γιατί είναι ένα δίκτυο που αφορούσε το NAS που είναι state of the art και επίσης είναι κατεχόμενο για mobile, οπότε είναι ένας από τους καλύτερους και κέρδιζε για την εφαρμογή που το θέλαμε.

Network Architecture Search

NASNet Mobile



Το NASNetMobile γράφει test set accuracy 62% στις 12 epochs και μετά φαίνεται να πέφτει.

Fig.3

Και δεν είναι μόνο το NASNetMobile στο οποίο παρατηρούμε αυτό. Η ερώτηση που μπορεί να δοθεί είναι η ερώτηση: πού είναι τα data που έχουμε δημιουργήσει είναι fake και προσομοιώνει από αυθεντικές εικόνες με τους περισσότερους. Αυτό σημαίνει πολύ στο validation acc. γιατί υπολογίζεται το δίκτυο σε βέλτη εικόνα που δεν έχει found, αλλά βλέπει unseen εικόνες που έχει found. Έτσι εφύγεται το ότι το validation acc γίνεται και αυτό το $\approx 100\%$. Τώρα όσον αφορά το test set accuracy, στο Fig.3 φαίνεται πως μετά από ένα σύνολο το δίκτυο αρχίζει να πέφτει. Επιπλέον είναι να επιβεβαιώσουμε ότι η ερώτηση του test set είναι ένας εικονικός που δεν έχει found no 76. Το δίκτυο, λοιπόν, μετά από ένα σύνολο φαίνεται να χάνει τα χαρακτηριστικά του καθαρού, γεγονός που δείχνει problema!

Προκρίνεται να γίνει έγκοψη η σύγκριση των εκπαιδευτικών δικτύων (και η εν δυνάμει 7
μελλοντική χρήση κινήσεων σε ensembling), αλλά αποθηκεύεται στον φάκελο checkpoint.

Ο φάκελος που ακολουθείται είναι ο εξής: ο κάθε φάκελος αποτελείται με βάση το
όνομα του δικτύου και μέσα περιέχει φάκελο που αποτελείται με βάση τον αριθμό των
κλάσεων που συζητήσαμε στην ζήτηση. Ο φάκελος αυτός (π.χ. 100-classes) περιέχει
μέσα το main.py σε .h5 αρχείο και ένα .log αρχείο που περιέχει το ιστορικό των
εκπαιδευτών.

Ακόμη, ο φάκελος code περιέχει όλα τα αρχεία κώδικα που έχουν γραφτεί - παλιά
και καινούρια.

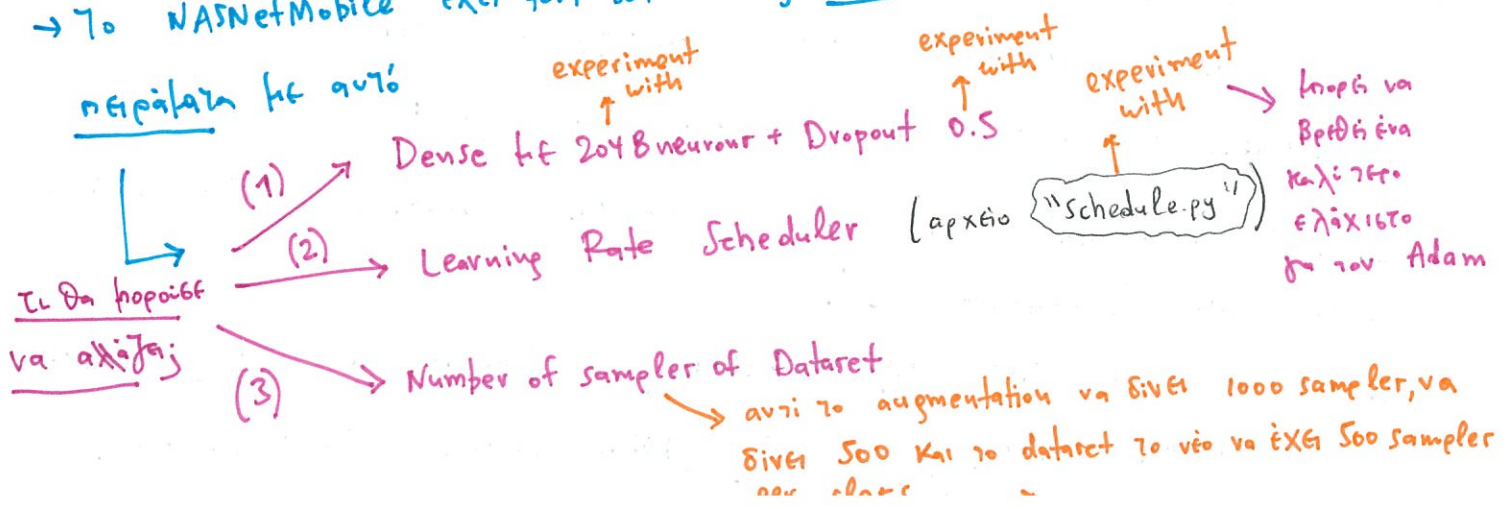
Το PlantNet είναι μια προεπιλεγμένη λειτουργία ενός κινητού custom νευρωνικού, το οποίο
μπορεί να αποδίδει ικανοποιητικά, βέβαια, και το NASNetMobile που παρουσιάστηκε παραπάνω
είναι δίκτυο κατάλληλο για χρήση σε κινητά. Στο φάκελο code υπάρχει το αρχείο

"tfLiteConverter.py", το οποίο μετατρέπει τα .h5 αρχεία σε tfLite αρχεία κατάλληλα
για χρήση σε κινητά.

ΣΥΝΟΨΗ + ΠΟΥ ΒΡΙΣΚΟΜΑΣΤΕ + ΕΠΟΜΕΝΑ ΒΗΜΑΤΑ

- Συνεχίζουμε περαιότερα με τις 100 κλάσεις με το augmented sampler 100 classes PlantCLEF2017val
- Αποδοκιμάζοντας, ένα μεγάλο δίκτυο μπορεί να φτάσει κοντά στο 80% test set accuracy 100 classes PlantCLEF2017Test
σε αυτό το 100 classes dataset. Αποφασίζουμε, όμως, να δοκιμάσουμε ένα πιο μικρό
για λόγους που αναφέραμε παραπάνω και εφόσον:

→ Το NASNetMobile έχει test set accuracy 62%, κρινεται σκόνη να συνεχιστούν
πειράματα με αυτό



Αυτό είναι πιθανό να οδηγήσει σε καλύτερη γενικευτική ικανότητα...

ΓΕΝΙΚΑ: και οι (3) αθροιστές ((1),(2),(3)) που προτείνονται έχουν δοκιμάσει στο paper και έχουν καταλήξει ότι αυτό είναι το καλύτερο,

ΑΛΛΑ χωρίς ακριβή ^{κρίση} περαιτέρω δοκιμές, βέβαια και προτείνονται!

- Το PlantNet (δουλειά το δικό μας custom CNN) θέλει πολλές δοκιμές για να βγάλει σε ακρίβεια του NASNetMobile, οπότε κρίνεται καλύτερο να ΜΗΝ εστιάζουμε σε αυτό (από βιβλία έχουν και κουνιστά/παράτοια όρα σε δίκτυα)
- Λόγω εναλλακτικής είναι να χρησιμοποιηθεί ένα κερτάρι δίκτυο, όπως αυτό που έχει βγάλει $\approx 80\%$ test accuracy και να γίνει προσπάθεια μέσω της λειτουργίας του .h5 αρχείου σε .tflite να περάσει ο όγκος του, ώστε να χωρέσει σε κινητό!

ΚΑΠΟΙΕΣ ΠΕΡΑΙΤΕΡΕ ΣΥΜΒΟΥΛΕΣ + TIPS

- Τα μικρά δίκτυα τρέχουν σχεδόν γρήγορα στο laptop του βασίλη (δουλειά τα 15 epochs του NASNetMobile θέλουν περίπου 4 ώρες. Τα κερτάρι δίκτυα είναι καλύτερα να τρέψουν στον υπολογιστή στην Ξάνθη

στον υπολογιστή αυτό έχει διατηρηθεί
ακριβώς ο ίδιος φορτωτής ονομαστικός ⁶¹⁰⁰⁷
6946005, αλλά έχουν τρέξει και τα keras CNN AnyDesk: 240656324
xanthi25410

- Στο φάκελο D:\advent\code υπάρχουν πολλές έτοιμες δίκτυα τα οποία

μπορούν να τρεχούν αλλά κάνουνται το αρχείο τους import στο "classification.py"

- ΥΠΟΨΙΑΖΟΜΑΙ πως θα είναι σχεδόν αδύνατο να βγάλουμε παραπάνω test set accuracy ή αν βγάλουμε θα βγάλουμε ελάχιστο. Στην περίπτωση που χρειάζομαστε παραπάνω πράξη να βρεθείτε τα εξής:

1. Το δικό μας dataset θα είναι πιο balanced και με καλύτερη (ελαφρύ) εικόνα, το οποίο θα συνηθίζεται ↗.
2. Στο paper του αυτό ΔΕΝ ΚΑΤΑΦΕΡΕ να ξεπεράσει το 65% validation set accuracy ΧΩΡΙΣ ENSEMBLE (βέβαια εκείνος είχε 10000 κλάσεις)

9
ὡς φαίνεται στο Submission Models Overview.xlsx. Ακόμα, δοκιμάσε
να χρησιμοποιήσεις τον αρχικοποιητή του τελικού του δικτύου βάσει αυτού
προσφοίμενα, και το οποίο μπορεί να δοκιμαστεί και σε συνδυασμό με
ENSEMBLE να αποδώσει. ΑΥΤΗ ίσως να είναι η καλύτερη λύση και
αυτή που θα αποδώσει το τελικό ΒΕΛΤΙΣΤΟ ΔΙΚΤΥΟ.