

ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ & ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ

ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ
Χειμερινό Εξάμηνο 2019-20

Προπαρασκευή 2ης Εργαστηριακής Άσκησης:
Αναγνώριση φωνής με Κρυφά Μαρκοβιανά Μοντέλα και Αναδρομικά Νευρωνικά Δίκτυα

ΠΕΡΙΓΡΑΦΗ

Σκοπός είναι η υλοποίηση ενός συστήματος επεξεργασίας και αναγνώρισης φωνής, με εφαρμογή σε αναγνώριση μεμονωμένων λέξεων. Το πρώτο μέρος αποσκοπεί στην εξαγωγή κατάλληλων ακουστικών χαρακτηριστικών από φωνητικά δεδομένα, χρησιμοποιώντας τα κατάλληλα πακέτα `python`, καθώς και η ανάλυση και απεικόνισή τους με σκοπό την κατανόηση και την εξαγωγή χρήσιμων πληροφοριών από αυτά. Τα εν λόγω χαρακτηριστικά είναι στην ουσία ένας αριθμός συντελεστών `cepstrum` που εξάγονται μετά από ανάλυση των σημάτων με μια ειδικά σχεδιασμένη συστοιχία φίλτρων (`filterbank`). Η συστοιχία αυτή είναι εμπνευσμένη από ψυχοακουστικές μελέτες.

Πιο συγκεκριμένα, το σύστημα που θα αναπτύξετε αφορά σε αναγνώριση μεμονωμένων ψηφίων (`isolated digits`) στα Αγγλικά. Τα δεδομένα που θα χρησιμοποιήσετε περιέχουν εκφωνήσεις 9 ψηφίων από 15 διαφορετικούς ομιλητές σε ξεχωριστά `.wav` αρχεία. Συνολικά θα βρείτε 133 αρχεία, αφού 2 εκφωνήσεις θεωρήθηκαν προβληματικές και δεν έχουν συμπεριληφθεί. Τα ονόματα των αρχείων (π.χ. `eight8.wav`) υποδηλώνουν τόσο το ψηφίο που εκφωνείται (π.χ. `eight`), όσο και τον ομιλητή (οι ομιλητές είναι αριθμημένοι από 1-15). Οι εκφωνήσεις έχουν ηχογραφηθεί με συχνότητα δειγματοληψίας ίση με $F_s = 16\text{kHz}$ και η διάρκειά τους διαφέρει.

ΒΙΒΛΙΟΘΗΚΕΣ PYTHON

- Διάβασμα αρχείων ήχου και εξαγωγή χαρακτηριστικών: `librosa`
- Αλγόριθμοι ταξινόμησης: `scikit-learn`
- Διαγράμματα: `matplotlib`, `seaborn` etc.
- Νευρωνικά δίκτυα: `pytorch`

ΕΠΙΠΛΕΟΝ ΛΟΓΙΣΜΙΚΟ

- Praat: <http://www.fon.hum.uva.nl/praat/>

ΕΚΤΕΛΕΣΗ

Κατεβάστε τα δεδομένα της προπαρασκευής από τις διευκρινίσεις του `mycourses`.

Βήμα 1

Ανάλυση αρχείων ήχου με το Praat (το οποίο πρέπει να εγκαταστήσετε από το παραπάνω link): Ανοίξτε τα αρχεία `onetwothree1.wav` και `onetwothree8.wav` με το πρόγραμμα Praat. Τα αρχεία αυτά περιέχουν την πρόταση “one two three” από τους ομιλητές 1 και 8, οι οποίοι είναι άντρας και γυναίκα αντίστοιχα. Παρατηρήστε τις κυματομορφές και τα `spectrograms` και έπειτα εξάγετε τη μέση τιμή του `pitch` στα φωνήεντα “α”, “ου”, “ι” για τα 3 ψηφία και για κάθε ομιλητή. Έπειτα, εξάγετε τα 3 πρώτα `formants` του κάθε φωνήεντος. Παρουσιάστε τα αποτελέσματα και γράψτε τις παρατηρήσεις σας.

Χρησιμοποιώντας Python 3, εκτελέστε τα παρακάτω βήματα:

Βήμα 2

Φτιάξτε μία συνάρτηση (data parser) που να διαβάζει όλα τα αρχεία ήχου που δίνονται μέσα στο φάκελο *digits/* και να επιστρέφει 3 λίστες Python, που να περιέχουν: Το wav που διαβάστηκε με librosa, τον αντίστοιχο ομιλητή και το ψηφίο.

Βήμα 3

Εξάγετε με το librosa τα Mel-Frequency Cepstral Coefficients (MFCCs) για κάθε αρχείο ήχου. Εξάγετε 13 χαρακτηριστικά ανά αρχείο. Χρησιμοποιήστε μήκος παραθύρου 25 ms και βήμα 10 ms. Επίσης, υπολογίστε και την πρώτη και δεύτερη τοπική παράγωγο των χαρακτηριστικών, τις λεγόμενες deltas και delta-deltas (hint: υπάρχει έτοιμη υλοποίηση στο librosa).

Βήμα 4

Αναπαραστήστε τα ιστογράμματα του 1ου και του 2ου MFCC των ψηφίων n1 και n2 για όλες τους τις εκφωνήσεις. Πόση απόκλιση υπάρχει?

Εξάγετε για 2 εκφωνήσεις των n1 και n2 από 2 διαφορετικούς ομιλητές τα Mel Filterbank Spectral Coefficients (MFSCs), δηλαδή τα χαρακτηριστικά που εξάγονται αφού εφαρμοστεί η συστοιχία φίλτρων της κλίμακας Mel πάνω στο φάσμα του σήματος φωνής αλλά χωρίς να εφαρμοστεί στο τέλος ο μετασχηματισμός DCT (εξάγετε και πάλι χαρακτηριστικά διάστασης 13). Αναπαραστήστε γραφικά τη συσχέτιση των MFSCs για την κάθε εκφωνήση. Σε ξεχωριστά διαγράμματα πραγματοποιήστε το ίδιο για τα MFCCs. Τι παρατηρείτε? Γιατί χρησιμοποιούμε τα MFCCs αντί των MFSCs?

Βήμα 5

Μια πρώτη προσέγγιση για την αναγνώριση των ψηφίων είναι η εξαγωγή ενός μοναδικού διανύσματος χαρακτηριστικών για κάθε εκφωνήση.

Ενώστε τα mfccs – deltas – delta-deltas και έπειτα για κάθε εκφωνήση δημιουργήστε ένα διάνυσμα παίρνοντας τη μέση τιμή και την τυπική απόκλιση κάθε χαρακτηριστικού για όλα τα παράθυρα της εκφωνήσης. Αναπαραστήστε με scatter plot τις 2 πρώτες διαστάσεις των διανυσμάτων αυτών, χρησιμοποιώντας διαφορετικό χρώμα και σύμβολο για κάθε ψηφίο. Σχολιάστε το διάγραμμα.

Βήμα 6

Μια καλή τακτική για απεικόνιση πολυδιάστατων διανυσμάτων είναι η μείωση των διαστάσεών τους με Principal Component Analysis (PCA). Μειώστε σε 2 τις διαστάσεις των διανυσμάτων του προηγούμενου βήματος με PCA και δημιουργήστε εκ νέου το scatter plot. Σχολιάστε και επαναλάβετε τη διαδικασία για 3 διαστάσεις και τρισδιάστατο scatter plot.

Τι ποσοστό της αρχικής διασποράς διατηρούν οι συνιστώσεις που προέκυψαν? Τι πληροφορία δίνουν αυτά τα νούμερα για τα principal components? Είναι επιτυχημένη η μείωση διαστάσεων?

Βήμα 7

Χωρίστε τα δεδομένα σε train-test με αναλογία 70%-30%. Ταξινομήστε με χρήση του Bayesian ταξινομητή της πρώτης εργαστηριακής άσκησης, καθώς και του Naive Bayes του scikit-learn. Χρησιμοποιήστε επίσης, άλλους 3 ταξινομητές της επιλογής σας. Αναφέρετε το ποσοστό επιτυχίας στο test set και συγκρίνετε τα αποτελέσματα. Σημείωση: Τα δεδομένα πριν την ταξινόμηση πρέπει να κανονικοποιηθούν.

(Bonus: Θα αυξηθεί το ποσοστό επιτυχίας αν προσθέσω επιπλέον ηχητικά χαρακτηριστικά στο διάνυσμά μου, όπως π.χ. zero-crossing rate? Χρησιμοποιήστε ελεύθερα τέτοια επιπλέον χαρακτηριστικά και εάν αυξηθεί το ποσοστό επιτυχίας αναφέρετε τη σχετική αύξηση, διαφορετικά αναφέρετε τους λόγους που απέτυχε η προσπάθειά σας)

Βήμα 8

Εξοικείωση με το PyTorch:

Δημιουργήστε ακολουθίες 10 σημείων ενός ημιτόνου και ενός συνημιτόνου με συχνότητα $f = 40$ Hz. Σκοπός

είναι η πρόβλεψη του συνημιτόνου με δεδομένη την ακολουθία του ημιτόνου. Επιλέξτε σταθερή και μικρή απόσταση ανάμεσα στα διαδοχικά σημεία.

Εκπαιδεύστε ένα Αναδρομικό Νευρωνικό Δίκτυο (Recurrent Neural Network – RNN), το οποίο θα δέχεται ως είσοδο τις ακολουθίες του ημιτόνου και θα πρέπει να προβλέπει τις αντίστοιχες ακολουθίες συνημιτόνου. Αντί για χρήση του απλού RNN μπορούν να χρησιμοποιηθούν και οι μονάδες LSTM και GRU (δώστε το λόγο που τις χρησιμοποιήσατε και γιατί είναι τόσο διαδεδομένες).

ΣΗΜΕΙΩΣΗ

Τα ψηφία n_1 και n_2 που αναφέρονται παραπάνω, είναι το προτελευταίο και το τελευταίο ψηφίο του Α.Μ. σας αντίστοιχα. Αν κάποιο είναι το 0, τότε ορίστε το ως το προηγούμενο ή επόμενο του άλλου ψηφίου, το οποίο δεν είναι μηδενικό.

ΠΑΡΑΔΟΤΕΑ

(1) Σύντομη αναφορά (σε pdf ή jupyter notebook) που θα περιγράφει τη διαδικασία που ακολουθήθηκε σε κάθε βήμα, καθώς και τα σχετικά αποτελέσματα. Τα αποτελέσματα πρέπει να συνοδεύονται και από ερμηνεία – σχολιασμό.

(2) Κώδικας Python (συνοδευόμενος από σύντομα σχόλια). Προσπαθήστε να κάνετε vectorized υλοποιήσεις.

Συγκεντρώστε τα (1) και (2) σε ένα .zip αρχείο το οποίο πρέπει να αποσταλεί μέσω του mycourses.ntua.gr πριν από τη διεξαγωγή του εργαστηρίου.

ΠΑΡΑΡΤΗΜΑ: Οδηγίες εγκατάστασης πακέτων Python

1η εναλλακτική: Χρήση του Virtual Machine και του υλικού από το εισαγωγικό εργαστήριο Python: <https://github.com/georgepar/python-lab>

2η εναλλακτική: Εγκατάσταση με miniconda: <https://conda.io/miniconda.html>

Ο τρόπος αυτός είναι εύκολος και για Linux και για Windows χρήστες.

- Εγκαταστήστε το miniconda, σε Python 3, και έπειτα ανοίξτε ένα terminal.
- Μπορείτε να εγκαταστήσετε τα πακέτα στο default περιβάλλον του miniconda, αλλιώς μπορείτε να δημιουργήσετε ένα καινούργιο περιβάλλον με την εντολή:

`conda create -n my_env pip`

Για να ενεργοποιήσετε το περιβάλλον αυτό η εντολή είναι:

`source activate my_env`

- Για χρήστες Windows εμφανίζεται ένα καινούργιο εκτελέσιμο πρόγραμμα μετά την εγκατάσταση εάν κάνετε search: Open Anaconda Prompt. Όταν ανοίξει το terminal ακολουθείτε την ίδια διαδικασία με τους χρήστες Linux.
- Για την εγκατάσταση των πακέτων που χρειάζονται για την άσκηση η εντολή είναι η εξής:

`pip install numpy matplotlib librosa scikit-learn torch torchvision`

- Για όσους επιθυμούν να δουλέψουν σε jupyter notebook θα πρέπει να τρέξουν επίσης:

`pip install jupyter`

και για να ανοίξουν τον notebook server:

`jupyter notebook`