

Ασκήσεις μελέτης της 23^{ης} διάλεξης

23.1. Θέλουμε να χρησιμοποιήσουμε μια τροποποιημένη μορφή του συνελικτικού νευρωνικού δικτύου της διαφάνειας 23 (LeNet), για να εντοπίζουμε τις συντεταγμένες (x, y) του κέντρου του κεφαλιού και των δύο ώμων σε εικόνες (ή video frames) που περιλαμβάνουν έναν μόνο άνθρωπο μπροστά από μια κονσόλα ηλεκτρονικών παιχνιδιών εφοδιασμένη με έγχρωμη κάμερα και κάμερα βάθους.¹ Η κάθε εικόνα έχει ανάλυση 256×256 και τέσσερα κανάλια (RGB και βάθος), δηλαδή είναι ένας τανυστής (tensor) τριών αξόνων, με σχήμα (shape) $(256, 256, 4)$. Όπως στο σχήμα της διαφάνειας 23, υπάρχουν δύο συνελικτικά στρώματα (convolutional layers) που παράγουν 6 και 16 χάρτες χαρακτηριστικών (feature maps) αντίστοιχα αλλά οι συνελίξεις χρησιμοποιούν πυρήνες (kernels) με παράθυρο 3×3 και είναι ευρείες (wide, same), δηλαδή χρησιμοποιούν padding και διατηρούν την ανάλυση της αρχικής εικόνας σε κάθε κανάλι (βλ. και διαφάνεια 10). Τα δύο στρώματα υποδειγματοληψίας (pooling) χρησιμοποιούν max-pooling με παράθυρο 4×4 και βήμα (stride) 4 και στους δύο άξονες. Τα δύο πρώτα (τα κρυφά) πυκνά (dense) στρώματα του τελικού MLP εξακολουθούν να έχουν 120 και 84 νευρώνες αντίστοιχα.

α) Πόσους πυρήνες θα χρησιμοποιεί το πρώτο συνελικτικό στρώμα και τι σχήμα θα έχει ο καθένας;

Απάντηση: Το πρώτο συνελικτικό στρώμα θα χρησιμοποιεί 6 πυρήνες, ώστε να προκύπτουν 6 χάρτες χαρακτηριστικών, όπως φαίνεται στο σχήμα της διαφάνειας 23. Ο κάθε πυρήνας θα έχει 4 φέτες (slices), αφού η είσοδος έχει τώρα 4 κανάλια. Γνωρίζουμε από την εκφώνηση ότι κάθε πυρήνας εφαρμόζει σε κάθε κανάλι της εισόδου παράθυρο 3×3 . Επομένως κάθε ένας από τους 6 πυρήνες θα είναι ένας τανυστής (tensor) τριών αξόνων, με σχήμα (shape) $(3, 3, 4)$.

β) Τι ανάλυση θα έχουν οι χάρτες χαρακτηριστικών που θα προκύπτουν από το πρώτο στρώμα max-pooling;

Απάντηση: Αφού τα συνελικτικά στρώματα που χρησιμοποιούμε διατηρούν την ανάλυση σε κάθε κανάλι, κάθε ένας από τους 6 χάρτες χαρακτηριστικών (κανάλια) που προκύπτουν από το πρώτο συνελικτικό στρώμα, δηλαδή κάθε κανάλι στην είσοδο του πρώτου στρώματος max-pooling θα εξακολουθεί να έχει ανάλυση 256×256 . Αφού κάθε στρώμα max-pooling χρησιμοποιεί παράθυρο 4×4 με βήμα (stride) 4 και στους δύο άξονες, ο κάθε ένας από τους 6 χάρτες που εξέρχονται από το πρώτο στρώμα max-pooling θα έχει ανάλυση $(256/4) \times (256/4)$, δηλαδή 64×64 .

γ) Πόσους πυρήνες θα χρησιμοποιεί το δεύτερο συνελικτικό στρώμα και τι σχήμα θα έχει ο καθένας;

Απάντηση: Το δεύτερο συνελικτικό στρώμα θα χρησιμοποιεί 16 πυρήνες, ώστε να προκύπτουν 16 χάρτες χαρακτηριστικών, όπως φαίνεται στο σχήμα της διαφάνειας 23. Ο κάθε πυρήνας θα έχει 6 φέτες (slices), αφού η είσοδος του συνελικτικού στρώματος (η έξοδος του πρώτου στρώματος max-pooling) έχει 6 κανάλια (χάρτες). Γνωρίζουμε από την εκφώνηση ότι κάθε πυρήνας εφαρμόζει σε κάθε κανάλι της εισόδου του παράθυρο 3×3 . Επομένως κάθε ένας από τους 16 πυρήνες θα είναι ένας τανυστής (tensor) τριών αξόνων, με σχήμα (shape) $(3, 3, 6)$.

¹ Υπάρχουν καλύτερα μοντέλα για αυτό το συγκεκριμένο πρόβλημα αλλά δεν περιλαμβάνονται στην ύλη αυτού του μαθήματος. Βλ. π.χ. την ενότητα 25.5 της 4^{ης} έκδοσης του βιβλίου «Τεχνητή Νοημοσύνη – Μια σύγχρονη προσέγγιση» των S. Russel και P. Norvig, Κλειδάριθμος, 2022.

δ) Τι ανάλυση θα έχουν οι χάρτες χαρακτηριστικών που θα προκύπτουν από το δεύτερο στρώμα max-pooling;

Απάντηση: Αφού τα συνελικτικά στρώματα που χρησιμοποιούμε διατηρούν την ανάλυση σε κάθε κανάλι, κάθε ένας από τους 16 χάρτες χαρακτηριστικών (κανάλια) που προκύπτουν από το δεύτερο συνελικτικό στρώμα, δηλαδή κάθε κανάλι στην είσοδο του δεύτερου στρώματος max-pooling θα εξακολουθεί να έχει ανάλυση 64×64 (όπως στην έξοδο του πρώτου στρώματος max-pooling). Αφού κάθε στρώμα max-pooling χρησιμοποιεί παράθυρο 4×4 με βήμα (stride) 4 και στους δύο άξονες, ο κάθε ένας από τους 16 χάρτες που εξέρχονται από το δεύτερο στρώμα max-pooling θα έχει ανάλυση $(64/4) \times (64/4)$, δηλαδή 16×16 .

ε) Πόσους νευρώνες θα έχει η είσοδος του τελικού MLP;

Απάντηση: Οι 16 χάρτες ανάλυσης 16×16 που εξέρχονται από το δεύτερο στρώμα max-pooling θα συνενώνονται σε ένα διάνυσμα $16 \times 16 \times 16 = 4096$ χαρακτηριστικών, που θα δίνεται ως είσοδος στο τελικό MLP (τρία πυκνά στρώματα) του σχήματος της διαφάνειας 23.

στ) Πόσους νευρώνες θα έχει το τελικό στρώμα εξόδου του MLP; Τι συνάρτηση ενεργοποίησης θα έχουν;

Απάντηση: Το στρώμα εξόδου του MLP θα έχει 6 νευρώνες, δύο για τις συντεταγμένες (x, y) του κεφαλιού και τέσσερις για τις συντεταγμένες των δύο ώμων. Οι νευρώνες αυτοί δεν θα έχουν συνάρτηση ενεργοποίησης, ώστε να μπορούν να παράγουν οποιονδήποτε πραγματικό αριθμό ο καθένας.

23.2. Μια εταιρεία κατασκευής οικιακών συσκευών ετοιμάζει έναν νέο τύπο (μοντέλο) φούρνου μικροκυμάτων που θα διαθέτει κάμερα. Η εταιρεία θέλει ο φούρνος να έχει τη δυνατότητα να αναγνωρίζει μέσω της κάμερας τον χρήστη που στέκεται μπροστά του, ώστε να προσαρμόζονται οι ρυθμίσεις του φούρνου στις προτιμήσεις του συγκεκριμένου χρήστη. Η εταιρεία σχεδιάζει να χρησιμοποιήσει ένα συνελικτικό νευρωνικό δίκτυο (CNN), το οποίο θα τροφοδοτείται με μια φωτογραφία του χρήστη που στέκεται μπροστά στη συσκευή. Το CNN θα έχει δέκα νευρώνες εξόδου, γιατί η εταιρεία θεωρεί ότι κάθε συσκευή του συγκεκριμένου τύπου θα χρησιμοποιείται σε ένα σπίτι ή γραφείο όπου οι χρήστες θα είναι το πολύ δέκα. Η εταιρεία διαθέτει 1.000 φωτογραφίες 50 ενδεικτικών χρηστών (20 από κάθε ενδεικτικό χρήστη) που έχουν τραβηγχεί με την κάμερα του νέου φούρνου. Κάθε μία από τις 1.000 φωτογραφίες είναι επισημειωμένη με τον κωδικό (id, 1–50) του αντίστοιχου ενδεικτικού χρήστη. Άλλα η εταιρεία δεν διαθέτει εκ των προτέρων φωτογραφίες όλων των χρηστών (σε κάθε σπίτι, γραφείο) που θα χρησιμοποιήσουν την κάθε μία συσκευή του συγκεκριμένου νέου τύπου. Όταν μία συσκευή του συγκεκριμένου τύπου εγκαθίσταται σε ένα σπίτι ή γραφείο, θα ζητείται από κάθε έναν από τους (το πολύ 10) χρήστες της να τραβήξει 5–10 φωτογραφίες του με την κάμερα της συσκευής, χρησιμοποιώντας ειδική επιλογή της διεπαφής χρήστη. Εξηγήστε πώς θα μπορούσε η εταιρεία να χρησιμοποιήσει τις 1.000 φωτογραφίες ενδεικτικών χρηστών που διαθέτει, καθώς και μια γενική συλλογή εκατομμυρίων επισημειωμένων εικόνων (π.χ. εικόνες ζώων, τοπίων κ.λπ., όπως στο ImageNet), ώστε να προ-εκπαιδεύσει (από το εργοστάσιο) το CNN του νέου τύπου φούρνου και να καταφέρει η κάθε συσκευή του νέου τύπου να αναγνωρίζει (με ελάχιστη πρόσθετη εκπαίδευση) τους συγκεκριμένους χρήστες της (σε συγκεκριμένο σπίτι ή γραφείο) έχοντας στη διάθεσή της μόνο 5–10 φωτογραφίες του καθενός.