

CONVOLUTIONAL NETWORKS

ARE ANY NETWORKS EMPLOYING CONVOLUTION INSTEAD THAN MATRIX MULTIPLICATION

CONVOLUTION: $S(t) = \int x(a)w(t-a)da$ $x = \text{INPUT}$, $w = \text{KERNEL}$, $S = \text{FEATURE MAP}$ **DISCRETE:** $S[t] = \sum_{a=-\infty}^{+\infty} x[a]w[t-a]$ **IN ML:** KERNEL IS LEARNABLE PARAMS AND A PDF, FEATURE

2D: $S[i,j] = \sum_m \sum_n x[m,n]k[i-m,j-n] = \sum_m \sum_n x[i-m,j-n]w[m,n] = \sum_m \sum_n x[i+m,j+n]k[m,n]$ **COMMUTATIVITY** **CROSS-CORRELATION** SAME AS CONVOLUTION BUT KERNEL NOT FLIPPED

REASONS FOR HAVING CONVOLUTION IN ANNS.

- SPARSE CONNECTIVITY | INTERACTIONS** KERNELS ARE $L \times L$ INPUT MATRIX. EVERY ORIGINAL INPUT INTERACT WITH SMALL SET OF HIDDEN (KERNEL) NODES, ITS WEIGHTS. **LOT LESS STORAGE SPACE** **FASTER PROCESSING** **UNITS IN DEEP LAYERS INTERACT INDIRECTLY WITH MOST INPUT STILL.**
- PARAMETER SHARING** SAME TIED WEIGHTS. KERNEL IS USED AT EVERY INPUT POSITION. MORE STORAGE CAPACITY REDUCTION. DRAMATICALLY MORE EFFICIENT
- INVARIANCE TO TRANSLATION:** IF INPUT CHANGES, OUTPUT CHANGES SAME WAY. $f(x)$ INVARIANT TO g IF $f(g(x)) = g(f(x))$ **BILLIONS** FOR SPATIAL **10,000** COMBINATIONS FOR 1000 X 1000 IMGS
- CONVOLUTIONAL LAYER:** CONVOLUTIONAL STAGE (AFFINE TRANSFORM) \rightarrow DETECTOR STAGE (NONLINEARITY) \rightarrow POOLING **OR** CONSIDER EACH STAGE A LAYER
- ALL STAGES HAVE SAME KERNEL TENSOR** **IE RECU**

POOLING

- REDUCES OUTPUT OF NONLINEARITIES WITH A SUMMARY STATISTIC OF NEARBY OUTPUTS. POOLING TYPES: MAX, AVG, L2 NORM, WEIGHTED AVG FROM CENTER (ON RECTANGULAR NEIGHBORHOOD)
- MAKES REPRESENTATION INVARIANT TO SMALL TRANSLATIONS** \rightarrow POOLED OUTPUTS WON'T CHANGE. WE CARE MORE ABOUT PRESENCE OF FEATURE THAN ITS LOCATION.
- IF WE POOL OVER OUTPUTS OF DIFFERENT KERNELS WE CAN MAKE INVARIANCE TO DIFFERENT TRANSFORMATION (ROTATED DIGITS EXAMPLE)**
- FILTERS WITH STRIDE > 1 RESULT IN DOWNSAMPLING** \rightarrow EVEN MORE REDUCTION IN COMPUTATIONAL AND STORAGE REQUIREMENTS
- MAY POOL STUFF DYNAMICALLY, FOR CURSING, OR USE THE POOLING STRUCTURE ITSELF**
- POOLING CAN BE TRICKY FOR NETS WITH TOP/DOWN INFO, OUTLIER MACHINE, AND AUTOENCODERS**
- CAN BE SEEN AS INFINITELY STRONG PRIOR OVER WEIGHTS FOR A FULLY CONNECTED NET** SPECIFYING ~~THE~~ PARAMS MUST BE TIED, WEIGHT ZERO EXCEPT FOR KERNEL REGION, ETC

CONVOLUTION FUNCTION VARIANTS

- IN PAVN WE DO MANY CONVOLUTIONS IN PARALLEL ON A LAYER.** DIFFERENT FEATURES, MULTIPLE CHANNELS (RGB) **IMPLEMENTATION WITH 4D KERNEL TENSOR:** SPATIAL X, Y, CHANNEL INDEX IN BATCH.
- IS COMMUTATIVE IF INPUT CHANNELS = OUTPUT CHANNELS** **STRIDE:** SKIP SOME POSITIONS TO SAVE RESOURCES, SAMPLE EVERY S PIXELS
- ZERO PADDING** TO AVOID LOSING MUCH INFO OVER MANY CONVOLUTIONAL LAYERS BECAUSE VALID CONVOLUTION $(M \times M) \times (N \times N) \rightarrow (M-N+1) \times (M-N+1)$ VERY AM. IF WIDE KERNELS \rightarrow USED TO EITHER MAKE OUTPUT NOT SHORER (SAME) OR TO MAKE EVERY PIXEL VISITED N TIMES IN EACH DIRECTION (FULL) $M+N-1 \times \text{OUTPUT}$ $M+N-1$
- LOCALLY CONNECTED LAYERS:** 6D TENSOR. DIFFERENT KERNEL PER PATCH. 'UNSHARED CONVOLUTION' USEFUL WHEN FEATURES ARE OF SMALL PATCH BUT DIFFERENT FEATURES AT DIFFERENT LOCATIONS. IE IN A FACE WE WANT LOOK FOR MOUTH AT THE TOP.
- CHANNEL-TO-CHANNEL CONNECT** EACH OUTPUT CHANNELS TO SUBSET OF INPUT CHANNELS. FURTHER EFFICIENTATION.
- TILED CONVOLUTION!** 'COMPROMISE' BETWEEN CONVOLUTIONAL AND LOCALLY CONNECTED LAYER. LEAVES A SET OF KERNELS WE ROTATE AS WE MOVE THROUGH SPACE **IF FILTER STACK IS SAME FEATURES TRANSFORMED AND THEN WE POOL** \rightarrow **ARBITRARILY TRANSFORMATION INVARIANT.**

HOW TO DO BACKPROP THROUGH CONVOLUTION.

- RECONSTRUCTION FUNCTION** MULTIPLY BY THE TRANSPOSE OF MATRIX DERIVED BY CONVOLUTION (ALSO IN AUTOENCODERS, RNN). CAREFUL BECAUSE NEEDS TO TAKE PADDING POLICY INTO ACCOUNT; ALSO STRIDE AND SIZE CONSERVING.
- CONVOLUTION** **GRADIENT TO WEIGHTS** **GRADIENT TO INPUT** IS ALL IT'S NEEDED TO TRAIN CONVNETS
- BIASES:** ONE PER CHANNEL, SHIMMED | ONE PER STACK | DIFFERENT LAYERS

STRUCTURED OUTPUTS

WHEN WE WANT TO OUTPUT A HIGH-DIM STRUCTURED OBJECT \rightarrow IE A PIXEL SEGMENTATION MAP. VECTOR OF SCORES (CATEGORY PROBS) PER PIXEL/PATCH

- REMOVE FULLY CONNECTED TOP LAYERS! GETS A SPATIALLY STRUCTURED OUTPUT, TRAIN WITH TARGET LABELS, APPLY SOFTMAX. LAST CONV LAYER MUST HAVE SAME DIMS AS NO. OUTPUT CATEGORIES
- \rightarrow WE MIGHT WANT TO REMOVE POOLING, OR SHARE KERNEL WEIGHTS ACROSS LAYERS (RESULTS IN SPECIAL RNN)
- \rightarrow THEN USE CONVNET OUTPUT TO TRAIN/FIT A CRF FOR INSTANCE (STRONG PRIOR NEIGHBORS HAVE SAME VALUE) RECURRENTS

DATA TYPES

| | SINGLE CHANNEL | MULTI CHANNEL | |
|----|---|-------------------------|--|
| 1D | AUDIO WAVEFORM | SKELETON ANIMATION DATA | • CONVNETS CAN BE USED FOR SOURCE W/ DIFFERENT SIZE BY MAKING IE POOLING PROPORTIONAL TO INPUT SIZE. |
| 2D | AUDIO IN FOURIER DOMAIN \rightarrow SPECTROGRAM | COLOR IMAGE DATA | |
| 3D | VOLUMETRIC DATA, CT SCAN | COLOR VIDEO DATA | |

CONVOLUTION ALGORITHMS

- CONVOLUTION IS MULTIPLICATION IN FOURIER DOMAIN. A LOT BETTER AND FASTER THEN DOING IT EXPLICITLY WITH OUTER PRODUCTS
- IF D-DIMENSIONAL KERNEL IS OUTER PRODUCT OF D VECTORS, KERNEL IS SEPARABLE \rightarrow NAIVE CONVOLUTION INEFFICIENT
- \rightarrow COMPOSE D 1-DIM CONVOLUTIONS WITH EACH VECTOR. $O(W \times D)$ VS $O(W^D)$ NOT EVERY CONVOLUTION IS LIKE THIS.

RANDOM OR UNSUPERVISED FEATURES

- IF NO SUPERVISED TRAINING FOR FEATURES \rightarrow RANDOM INITIALIZATION, OR LEARN UNSUPERVISED. THEN USE THEM (INITIAL) FOR LAST CLASSIFIER LAYER(S)
- RANDOM WORKS WELL \rightarrow NATURALLY COMPOSE TOWARDS FREQUENCY SELECTIVITY AND INVARIANCE.
- USE GREEDY UNSUPERVISED PRETRAINING \rightarrow LAYER PER LAYER, LIKE MIF \rightarrow CONVOLUTIONAL DEEP BELIEF NET
- TRAIN DENSE UNSUPERVISED MODEL ON IMAGE PATCH, THEN USE THEM AS KERNEL FOR CONVOLUTIONAL LAYER
- \rightarrow WE CAN TRAIN CONVNET WITHOUT EVEN USING CONVOLUTION AT TRAINING \rightarrow ONLY IN INFERENCE!

NEUROSCIENTIFIC BASIS FOR CONVNETS

FINDINGS FROM VISUAL SYSTEMS NEUROSCIENCE (WIDEL-WASHY). EATS!! RETINA \rightarrow OPTIC NERVE \rightarrow LATERAL GENICULATE NUCLEUS \rightarrow V1 AREA \rightarrow V2 \rightarrow V4 \rightarrow IT
INFEROTEMPORAL CORTEX

A CONVNET IS INSPIRED BY STUFF IN V1: 2D SPATIAL MAP OF RETINA, SIMPLE DETECTORS RECEPTIVE FIELDS, COMPLEX SHIFT-INVARIANT DETECTORS, POOLING

- \rightarrow GRANDMOTHER CELLS IN MEDIAL TEMPORAL LOBE FOR CONCEPTS. FIRE REGARDLESS OF SENSORY MODALITY. • IT: CLOSEST ANALOG TO CONVNET VIST WAY OF FEATURES
- \rightarrow V1 CELLS EXHIBIT RESPONSE FUNCTIONS SIMILAR TO GABOR FILTERS: DESCRIBE 2D WIGTHS AT IMAGE
- MANY ML ALSO LEARN GABOR-LIKE EDGE DETECTION PATTERNS

$$W = a \exp(-p_x x'^2 - p_y y'^2) \cos(fx' + \phi)$$

$$x' = (x - x_0) \cos(\theta) + (y - y_0) \sin(\theta)$$

$$y' = -(x - x_0) \sin(\theta) + (y - y_0) \cos(\theta)$$

PARAMS CONTROL TUNING, ROTATIONS, FREQUENCY RESPONSE

CONVNET HISTORY

CHECK READING SYSTEM (LEARN). OCR HANDWRITING BY MICROSOFT. FIRST DEEP MODEL TRAINED WITH BACKPROP. MICROSOFT DEEP LEARNING TODAY UP TODAY
THEIR EARLY SUCCESS PROBABLY DUE TO THEIR LARGE COMPUTATIONAL EFFICIENCY WRT FULLY CONNECTED NETS.