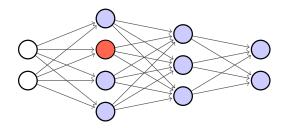
EPFL

Image Processing I

Chapter 6 Convolutional neural networks

Prof. Michael Unser, LIB



December 2022

OUTLINE

Introduction

- The (deep) learning (r)evolution in signal processing
- Artificial neurons
- Neural networks architectures for image processing

Basic Components of CNNs

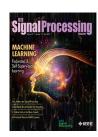
- Operator-based formalism
- Composition properties
- Pooling
- Continuity properties

CNNs in Practice

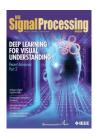
- Deep learning pipeline
- Denoising
- Segmentation

The (deep) learning (r)evolution in image processing

Special issues

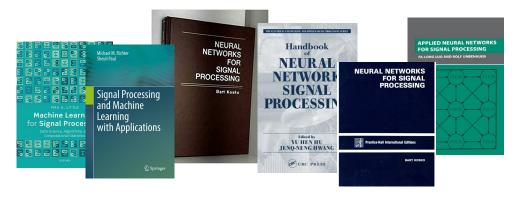


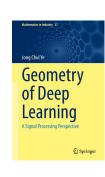






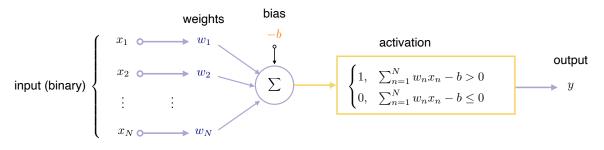
Flurry of new textbooks on neural networks

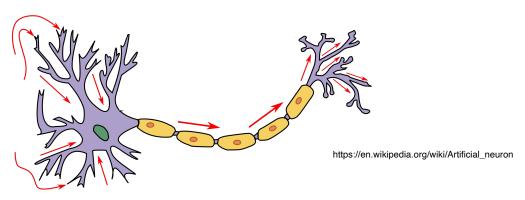




6-3

Formal model of neuron (McCulloch & Pitt)





Artificial neurons



Definition: An artificial neuron with weights $\mathbf{w}=(w_1,\ldots,w_N)\in\mathbb{R}^N$, bias $b\in\mathbb{R}$ and activation function $\sigma:\mathbb{R}\to\mathbb{R}$ is defined as the function $f:\mathbb{R}^N\to\mathbb{R}$

$$f(\boldsymbol{x}) = \sigma\left(\mathbf{w}^{\mathsf{T}}\boldsymbol{x} - b\right) = \sigma\left(\sum_{n=1}^{N} w_n x_n - b\right).$$

Examples of activation functions

- $\blacksquare \ \, \text{Threshold Logic Unit (Heaviside):} \quad \text{TLU}(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \, \text{(McCullogh \& Pitt 1943; Rosenblatt 1957)}$
- $\qquad \text{Sigmoid function:} \quad \sigma(x) = \frac{1}{1 + \mathrm{e}^{-x}}$

(Rumelhart 1986, ...)



- Rectified Linear Unit: $ReLU(x) = x_+ = max(0, x)$
- And variants





6-5

Neural network architectures for image processing

Neural networks are constructed from the composition of basic modules that can be chained at will.

Convolutional neural networks (CNN), in particular, are inspired by the structure of the primary visual cortex.

They have an architecture that is well suited for image processing.

Basic modules

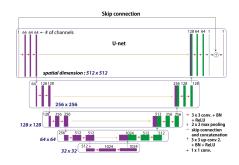
- Multi-channel convolution operators (filters)
- Pointwise nonlinearities
- Pooling: linear combination, flattening, sub-sampling, ...

Some modules—in particular, the filters—are adjustable.

The parameters of the CNN (weights) are set during the training procedure.



- Requires a comprehensive collection of reference input-output pairs.
 The larger the training set, the better!
- Formulated as a large-scale optimization problem
- Solved using some form of stochastic gradient algorithm (ADAM)
- Requires a lot of computational ressources (GPU)



Basic Components of CNNs

- Operator-based formalism
- Composition
- Pooling
- Continuity and stability estimates

Unser: Image processing 6-7

Operator-based formalism

Generic operator $T: \mathcal{X} \to \mathcal{Y}$ where \mathcal{X} and \mathcal{Y} are complete normed vector spaces.

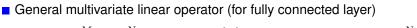
$$\forall x \in \mathcal{X}: \quad y = \mathrm{T}\{x\} \in \mathcal{Y}$$

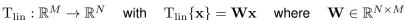
■ General multivariate nonlinear operator (for image patches)

$$\mathbf{T}:\mathbb{R}^M\to\mathbb{R}^N$$

■ Pointwise nonlinearity = activation function of neuron

$$T_{acti}: \mathbb{R} \to \mathbb{R}$$
 with $T_{acti}\{x\} = \sigma(x+b), b \in \mathbb{R}$





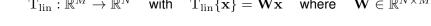




Image-to-image operator

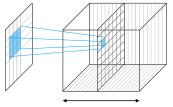
$$T: \ell_2(\mathbb{Z}^d) \to \ell_2(\mathbb{Z}^d)$$

Most operators of interest are shift-invariant

Example:
$$T_{LSI}\{f\} = h * f$$
 with $h = T_{LSI}\{\delta[\cdot]\}$ (discrete convolution)

Convolutional layer

■ Patch extraction operator: $\ell_2(\mathbb{Z}^d) \times \mathbb{Z}^d \to \mathbb{R}^M$ with M = #W $\operatorname{Vect}(f[\cdot], \mathbf{k}) = \mathbf{f}_W[\mathbf{k}] = (f[\mathbf{k} - \mathbf{k}_0])_{\mathbf{k}_0 \in W}$



■ Convolution layer : $\ell_2(\mathbb{Z}^d) \to \ell_2^N(\mathbb{Z}^d)$ with N channels = feature maps

Shared operator $T_{\mathrm{patch}}: \mathbb{R}^M \to \mathbb{R}^N$

T_{patch}
$$(m{f}_W[m{k}]) = m{\sigma}ig(\mathbf{W} m{f}_W[m{k}] ig)$$
 where $\mathbf{W} = egin{pmatrix} \mathbf{w}_1^\mathsf{T} \\ \vdots \\ \mathbf{w}_N^\mathsf{T} \end{pmatrix}$ and $m{\sigma} = egin{pmatrix} \sigma_1 : \mathbb{R} \to \mathbb{R} \\ \vdots \\ \sigma_N : \mathbb{R} \to \mathbb{R} \end{pmatrix}$

Implementation of N channel convolution layer:

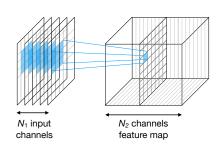
$$\begin{pmatrix} \sigma_1\Big((h_1*f)[\boldsymbol{k}]\Big) \\ \vdots \\ \sigma_N\Big((h_N*f)[\boldsymbol{k}]\Big) \end{pmatrix} = \mathrm{T}_{\mathrm{patch}}(\boldsymbol{f}_W[\boldsymbol{k}]) = \boldsymbol{\sigma}\Big(\mathbf{W}\boldsymbol{f}_W[\boldsymbol{k}]\Big)$$

6-9

Vector-valued convolutional layer (tensor)

Input feature map: $f[\cdot] = (f_1[\cdot], \dots, f_{N_1}[\cdot])$

■ Tensor patch extraction: $\ell_2^{N_1}(\mathbb{Z}^d) \times \mathbb{Z}^d \to \mathbb{R}^{M \times N_1}$ with M = #WTensor($\mathbf{f}[\cdot], \mathbf{k}$) = $\mathbf{F}[\mathbf{k}] = (f_i[\mathbf{k} - \mathbf{k}_0])_{\mathbf{k}_0 \in W, i \in \{1, ..., N_1\}}$



■ Vector-valued convolution layer : $\ell_2^{N_1}(\mathbb{Z}^d) \to \ell_2^{N_2}(\mathbb{Z}^d)$

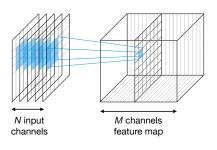
$$\text{Shared operator } \mathbf{T}_{\text{tensor}}: \mathbb{R}^{M \times N_1} \to \mathbb{R}^{N_2} \qquad \begin{array}{c} \text{convolution tensors} & \text{pointwise nonlinearities} \\ \mathbf{T}_{\text{tensor}}(\boldsymbol{F}[\boldsymbol{k}]) = \boldsymbol{\sigma} \big(\mathbf{W} \boldsymbol{F}[\boldsymbol{k}] \big) & \text{where} & \mathbf{W} = \begin{pmatrix} \mathbf{w}_1^\mathsf{T} \\ \vdots \\ \mathbf{w}_{N_2}^\mathsf{T} \end{pmatrix} \text{ and } \boldsymbol{\sigma} = \begin{pmatrix} \sigma_1 : \mathbb{R} \to \mathbb{R} \\ \vdots \\ \sigma_{N_2} : \mathbb{R} \to \mathbb{R} \end{pmatrix}$$

$$= m{\sigma} \Big((\mathbf{H} * m{f}) [m{k}] \Big) \quad ext{with} \quad \mathbf{H}[\cdot] = N_2 imes N_1 ext{ array of filters}$$

Convolution of vector-valued images

Input feature map:
$$m{f}[\cdot] = egin{pmatrix} f_1[\cdot] \\ \vdots \\ f_N[\cdot] \end{pmatrix}$$

$$\text{Matrix-valued filter:} \quad \mathbf{H}[\cdot] = \begin{pmatrix} h_{1,1}[\cdot] & h_{1,2}[\cdot] & \dots & h_{1,N}[\cdot] \\ \vdots & \vdots & & \vdots \\ h_{M,1}[\cdot] & h_{1,2}[\cdot] & \dots & h_{M,N} \end{pmatrix}$$



■ Vector-valued filterbank: $\ell_2^N(\mathbb{Z}^d) \to \ell_2^M(\mathbb{Z}^d)$

$$(\mathbf{H} * \mathbf{f})[\cdot] = \begin{pmatrix} (h_{1,1} * f_1)[\cdot] + (h_{1,2} * f_2)[\cdot] + \dots + (h_{1,N} * f_N)[\cdot] \\ \vdots \\ (h_{M,1} * f_1)[\cdot] + (h_{M,2} * f_2)[\cdot] + \dots + (h_{M,N} * f_N)[\cdot] \end{pmatrix}$$

6-11

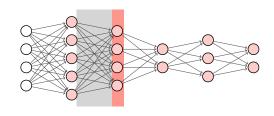
Composition

Series of (nonlinear) operators $T_\ell: \mathcal{X}_\ell \to \mathcal{Y}_\ell$ where \mathcal{X}_ℓ and \mathcal{Y}_ℓ are complete normed vector spaces.

Hypothesis: T_ℓ and $T_{\ell-1}$ have compatible domain and range; i.e., $\mathcal{X}_\ell = \mathcal{Y}_{\ell-1}$

 \blacksquare Composition of $T_1:\mathcal{X}_1\to\mathcal{X}_2$ and $T_2:\mathcal{X}_2\to\mathcal{X}_3$

$$\mathrm{T}_2\circ\mathrm{T}_1:\mathcal{X}_1\to\mathcal{X}_3\quad\text{with}\quad\mathrm{T}_2\circ\mathrm{T}_1\{x\}=\mathrm{T}_2\big\{\mathrm{T}_1\{x\}\big\}$$



Deep neural network

$$\mathbf{f}_{\mathrm{deep}}(\boldsymbol{x}) = \left(\boldsymbol{\sigma}_{L} \circ \mathbf{A}_{L} \circ \boldsymbol{\sigma}_{L-1} \circ \cdots \circ \boldsymbol{\sigma}_{2} \circ \mathbf{A}_{2} \circ \boldsymbol{\sigma}_{1} \circ \mathbf{A}_{1}\right)(\boldsymbol{x})$$

- Pointwise nonlinearities $\sigma_\ell: \mathbb{R}^{N_\ell} \to \mathbb{R}^{N_\ell}$ $\sigma_\ell(\boldsymbol{x}) = \left(\sigma(x_1), \dots, \sigma(x_{N_\ell})\right) \quad \text{with common activation fonction } \sigma: \mathbb{R} \to \mathbb{R}$

Composition: Properties

Preservation of linearity (affiness)

$$A_1: \mathbf{x} \mapsto \mathbf{W}_1\mathbf{x} + \mathbf{b}_1 \text{ and } A_2: \mathbf{y} \mapsto \mathbf{W}_2\mathbf{y} + \mathbf{b}_2 \quad \Rightarrow \quad A_2 \circ A_1: \mathbf{x} \mapsto (\mathbf{W}_2\mathbf{W}_1)\mathbf{x} + (\mathbf{b}_2 + \mathbf{W}_2\mathbf{b}_1)$$

Preservation of convolutional structure

 T_1 and T_2 are LSI with impulse responses $h_1, h_2 \in \ell_1(\mathbb{Z}^d)$

$$\Rightarrow$$
 $T_2 \circ T_1 : f \mapsto h * f$ with $h = h_2 * h_1 \in \ell_1(\mathbb{Z}^d)$

Application: Construction of larger receptive fields

Preservation of continuity

The composition of two continuous functions is continuous

6-13

Pooling







- Down-sampling: $\downarrow_m \{f\}[k] = f[mk]$
- lacksquare Max pooling: $\mathbb{R}^N o \mathbb{R}$ $\boldsymbol{u} \mapsto \max(\boldsymbol{u}) = \max(u_1, \dots, u_N)$
- Softmax: $\mathbb{R}^N \to \mathbb{R}^N$ (transforms output of CNN in "probabilities")

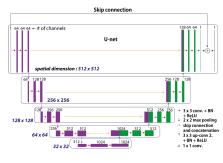
$$\mathbf{u} \mapsto \left(p_k = \frac{\exp(u_k)}{\sum_{n=1}^N \exp(u_n)} \right)$$

Up-sampling

$$\uparrow_m\{f\}[\pmb{k}] = \begin{cases} f[\pmb{n}], & \pmb{k} = m\pmb{n} \\ 0, & \text{otherwise} \end{cases} \qquad \uparrow_m\{f\}[\pmb{k}] = f[\pmb{k}/m]$$

Up-sampling with repetition

$$\uparrow_m\{f\}[\mathbf{k}] = f[\mathbf{k}/m]$$



U-net

Continuity requirements

Generic operator $T: \mathcal{X} \to \mathcal{Y}$ where $(\mathcal{X}, \|\cdot\|_{\mathcal{X}})$ and $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}})$ are complete normed vector spaces.

Definition

The operator $T: \mathcal{X} \to \mathcal{Y}$ is **continuous** if $\lim_i T\{x_i\} = T\{\lim_i x_i\} = T\{x\}$ for any sequence $(x_i)_{i \in \mathbb{N}}$ that is converging in \mathcal{X} with $\lim_i x_i = x$

Definition

The operator $T: \mathcal{X} \to \mathcal{Y}$ is **Lipschitz continuous** if there exists a constant L > 0 such that $\|T\{x_1\} - T\{x_2\}\|_{\mathcal{Y}} \le L\|x_1 - x_2\|_{\mathcal{X}}$ for any $x_1, x_2 \in \mathcal{X}$.

Lipschitz constant

 $\operatorname{Lip}(T) = L$ where L is the smallest constant such that the Lipschitz inequality holds.

To avoid instabilities, all modules of a CNN should be Lipschitz continuous.

This implies that they are a.e. differentiable, which is desirable for training with backpropagation.

Counterexample: TLU networks are discontinuous, and therefore very hard to train unless the architecture is shallow (perceptron).

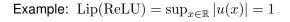
6-15

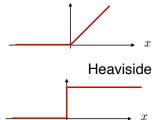
Lipschitz constant of primary modules

■ Pointwise nonlinearity

 $\sigma: \mathbb{R} \to \mathbb{R}$ where σ is differentiable

$$\operatorname{Lip}(\sigma) = \sup_{x \in \mathbb{R}} \left| rac{\mathrm{d}\sigma(x)}{\mathrm{d}x} \right| = \|\sigma'\|_{L_{\infty}}$$
 (cf. Mean Value Theorem)





LSI operator (convolution channel)

$$\mathrm{T}_{\mathrm{LSI}}(f) = h * f \quad ext{ where } h \in \ell_1(\mathbb{Z}^d)$$

$$\operatorname{Lip}(T_{\mathrm{LSI}}) = H_{\max} = \sup_{\boldsymbol{\omega} \in [0,\pi]^d} \left| H(e^{j\boldsymbol{\omega}}) \right| \le \|h\|_{\ell_1}$$

Justification (Parseval)

$$\begin{split} \|h*f - h*g\|_{\ell_2}^2 &= \|h*(f-g)\|_{\ell_2}^2 = \frac{1}{(2\pi)^d} \int_{[0,\pi]^d} |H(\mathrm{e}^{\mathrm{j}\boldsymbol{\omega}})|^2 \left| F(\mathrm{e}^{\mathrm{j}\boldsymbol{\omega}}) - G(\mathrm{e}^{\mathrm{j}\boldsymbol{\omega}}) \right|^2 \mathrm{d}\boldsymbol{\omega} \\ &\leq \frac{H_{\max}^2}{(2\pi)^d} \int_{[0,\pi]^d} \left| F(\mathrm{e}^{\mathrm{j}\boldsymbol{\omega}}) - G(\mathrm{e}^{\mathrm{j}\boldsymbol{\omega}}) \right|^2 \mathrm{d}\boldsymbol{\omega} = H_{\max}^2 \|f - g\|_{\ell_2}^2 \end{split}$$

Lipschitz constant of primary modules (cont'd)

Linear (resp. affine) transform

$$\begin{aligned} \mathrm{T_{lin}}: \mathbb{R}^M \to \mathbb{R}^N & \text{ with } & \mathbf{x} \mapsto \mathbf{A}\mathbf{x} \text{ (linear)} \\ & \text{or } & \mathbf{x} \mapsto \mathbf{A}\mathbf{x} + \mathbf{b} \text{ (affine)} & \text{where } \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M \end{aligned}$$

$$\operatorname{Lip}(T_{\operatorname{lin}}) = \sup_{\|\mathbf{x}\|_2 \leq 1} \|\mathbf{A}\mathbf{x}\|_2 = \rho(\mathbf{A}) \quad \text{(spectral norm = largest singular value of } \mathbf{A})$$

Imposing Lip-1 layer by spectral normalization

$$T_{normal}\{\mathbf{x}\} = \frac{1}{\rho(\mathbf{A})}\mathbf{A}\mathbf{x}$$

Estimation by power method: For $k = 0, \dots, K$

$$\mathbf{u}_k = \frac{1}{\|\mathbf{A}^\mathsf{T}(\mathbf{A}\mathbf{u}_k)\|_2} \mathbf{A}^\mathsf{T}(\mathbf{A}\mathbf{u}_k)$$

Upon convergence, $\mathbf{u} = \lim_k \mathbf{u}_k$ is the dominant eigenvector of $\mathbf{A}^\mathsf{T} \mathbf{A}$.

Finally,
$$\rho(\mathbf{A}) = \sqrt{(\mathbf{A}\mathbf{u})^{\mathsf{T}}(\mathbf{A}\mathbf{u})}$$

6-17

Stability and combination of modules

Composition

$$T_1$$
 T_2

$$\operatorname{Lip}(\mathbf{T}_1) = L_1 \& \operatorname{Lip}(\mathbf{T}_2) = L_2 \qquad \Rightarrow \qquad \operatorname{Lip}(\mathbf{T}_2 \circ \mathbf{T}_1) \leq L_2 L_1$$

$$\forall f,g \in \mathcal{X}_1: \left\| \mathbf{T}_2 \circ \mathbf{T}_1\{f\} - \mathbf{T}_2 \circ \mathbf{T}_1\{g\} \right\|_{\mathcal{Y}_2} \leq L_2 \left\| \mathbf{T}_1\{f\} - \mathbf{T}_1\{g\} \right\|_{\mathcal{Y}_1 = \mathcal{X}_2} \leq L_2 L_1 \|f - g\|_{\mathcal{X}_1} + L_2 \|f\|_{\mathcal{Y}_2} \leq L_2 \|f\|_{\mathcal{$$

Deep neural network with Lip -1 activations (e.g., ReLU) :

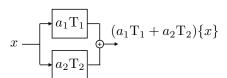
$$\mathbf{f}_{\mathrm{deep}}(\boldsymbol{x}) = \left(\boldsymbol{\sigma}_{L} \circ \mathbf{A}_{L} \circ \boldsymbol{\sigma}_{L-1} \circ \dots \circ \boldsymbol{\sigma}_{2} \circ \mathbf{A}_{2} \circ \boldsymbol{\sigma}_{1} \circ \mathbf{A}_{1}\right)(\boldsymbol{x})$$

ep neural network with
$$\operatorname{Lip}$$
-1 activations (e.g., ReLU) :
$$\mathbf{f}_{\operatorname{deep}}(\boldsymbol{x}) = (\boldsymbol{\sigma}_L \circ \operatorname{A}_L \circ \boldsymbol{\sigma}_{L-1} \circ \cdots \circ \boldsymbol{\sigma}_2 \circ \operatorname{A}_2 \circ \boldsymbol{\sigma}_1 \circ \operatorname{A}_1) (\boldsymbol{x}) \qquad \Rightarrow \qquad \operatorname{Lip}(\mathbf{f}_{\operatorname{deep}}) = \prod_{\ell=1}^L \left(\overbrace{\operatorname{Lip}(\boldsymbol{\sigma}_\ell)}^1 \operatorname{Lip}(\operatorname{A}_\ell) \right)$$

→ Deeper networks tend to be less stable

Linear combination

$$\operatorname{Lip} \big(\sum_{i=1}^I a_i \mathrm{T}_i \big) \leq \sum_{i=1}^I |a_i| \operatorname{Lip} (\mathrm{T}_i) \qquad \text{(by triangle inequality)}$$



Parallel feature maps

$$\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 \\ \vdots \\ \mathbf{T}_N \end{pmatrix} : \ell_2(\mathbb{Z}^d) \to \ell_2^N(\mathbb{Z}^d) \\ \operatorname{Lip}(\mathbf{T}) \le \sqrt{L_1^2 + \dots L_N^2} \le L_1 + \dots + L_N \qquad \qquad x \longrightarrow \begin{bmatrix} \mathbf{T}_1 \\ \vdots \\ \mathbf{T}_N \end{bmatrix} \xrightarrow{\mathbf{T}_1} \begin{pmatrix} \mathbf{T}_1\{x\} \\ \vdots \\ \mathbf{T}_N\{x\} \end{pmatrix}$$

$$x \longrightarrow \begin{bmatrix} T_1 \\ \vdots \\ T_N \end{bmatrix} \longrightarrow \begin{bmatrix} T_1\{x\} \\ \vdots \\ T_N\{x\} \end{bmatrix}$$

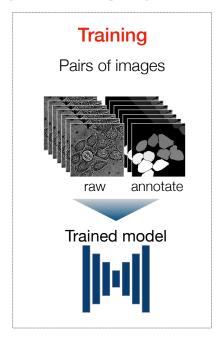
6-18

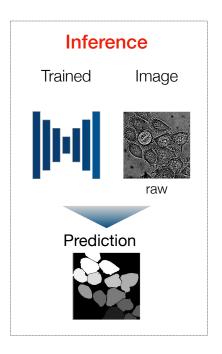
CNNs in Practice

- Deep learning pipeline
- Denoising
- Segmentation

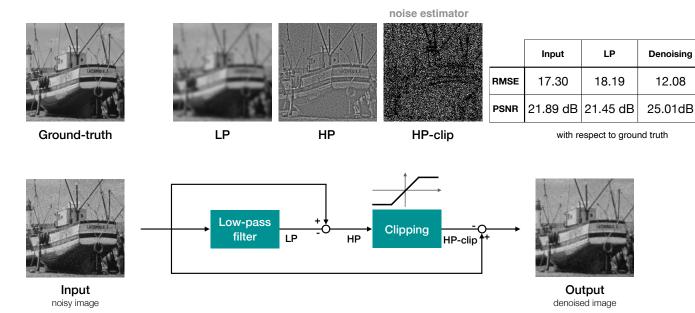
Unser: Image processing 6-19

Deep Learning Pipeline





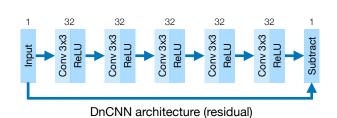
Denoising: "Handcrafted" ancestor of Resnet



6-21

Deep CNN for residual image denoisng







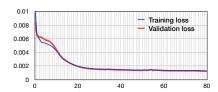
Datasets

• 80 natural graylevel images 384×384, artificially degraded by adding noise



Training CNN

- 20 images (25%) in validation set
- 37'601 parameters (weights and biais)
- Loss function is the MSE
- 80 epochs, batch_size = 16, lr = 0.001
- Global normalization 0/1



Layer	Shape	Parameters
Input	384 x 384 x 1	0
Conv2D 3x3 + b, ReLU	384 x 384 x 32	320
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9'248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9'248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9'248
Conv2D 3x3 + b, ReLU	384 x 384 x 32	9'248
Conv2D 3x3 + b, ReLU	384 x 384 x 1	289
Subtract	384 x 384 x 1	0
Total		37'601

Parameters

320

9248

9248

9248

9248

Shape

384 x 384 x 1

384 x 384 x 32

Layer

Input

Conv2D 3x3 + b, ReLU

Total

Ground-truth

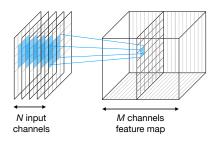
DnCNN architecture

First conv layer: N channels (feature maps)

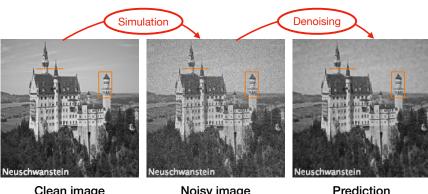
Number of learnable parameters: $(3 \times 3) + 1$ (bias) per output channel

Internal conv layers:

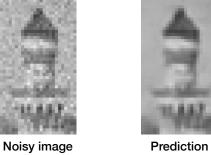
Clean image

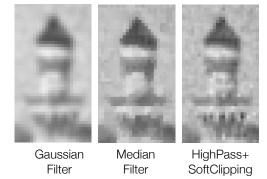


Number of learnable parameters: $(3 \times 3) \times 32 + 1$ (bias) per output channel







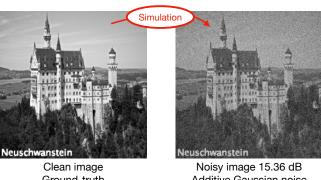


Gaussian noise $\sigma = 20$

Noisv

32

6-23



Ground-truth



Additive Gaussian noise



High-pass filter Soft Clipping 19.97 dB



Resnet (100 epochs) 3 layers, 16 channels 20.91 dB



Gaussian filter $\sigma = 1$ 18.23 dB



Median filter radius = 3 18.45 dB

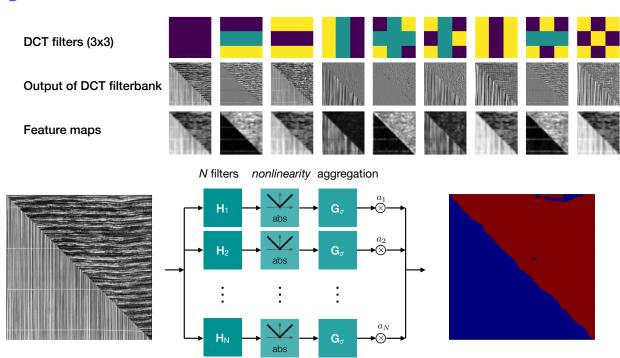


Resnet (200 epochs) 5 layers, 32 channels 21.01 dB



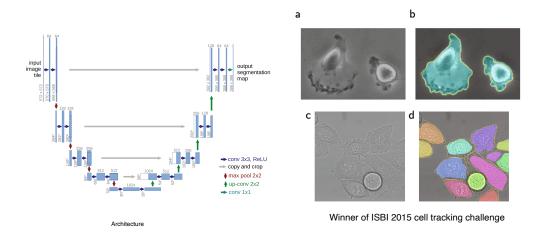
Unet (100 epochs) 3 Pooling steps, 32 channels 20.99 dB

Segmentation: "Handcrafted" texture discriminator



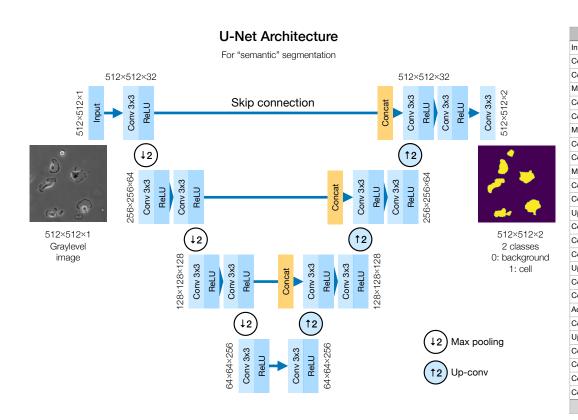
Popular CNN architecture for image segmentation

U-net introduced by Ronneberger in 2015 for biomedical image segmentation

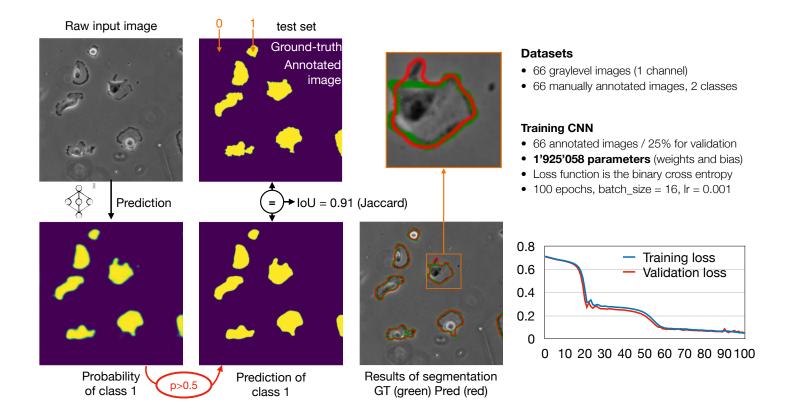


1. https://arxiv.org/abs/1505.04597

6-27



Layer	Shape	Params
nput	512x512x1	0
Conv2D 3x3	512x512x32	320
Conv2D 3x3, Activation	512x512x32	9248
/laxPooling2D	256x256x32	0
Conv2D 3x3, Activation	256x256x64	18496
Conv2D 3x3, Activation	256x256x64	36928
MaxPooling2D	128x128x64	0
Conv2D 3x3, Activation	128x128x128	73856
Conv2D 3x3, Activation	128x128x128	147584
MaxPooling2D	64x64x128	0
Conv2D 3x3, Activation	64x64x256	295168
Conv2D 3x3, Activation	64x64x256	590080
Jp-conv	128x128x128	131200
Concatenate	128x128x256	0
Conv2D 3x3, Activation	128x128x128	295040
Conv2D 3x3, Activation	128x128x128	147584
Jp-conv	256x256x64	32832
Concatenate	256x256x12	0
Conv2D 3x3, Activation	256x256x64	73792
Activation	256x256x64	0
Conv2D 3x3, Activation	256x256x64	36928
Jp-conv	512x512x32	8224
Concatenate	512x512x64	0
Conv2D 3x3, Activation	512x512x32	18464
Conv2D 3x3, Activation	512x512x32	9248
Conv2D 3x3	512x512x2	66
Total		1'925'058



CONCLUSION: The Pros and Cons of CNNs

- Advantages of (deep) CNNs
 - Are sufficiently flexible to implement most image-processing tasks
 - Can benefit from hardware acceleration (GPU)
 - Performance of deep CNNs is often spectacular
 - Once trained, they are very fast to deploy (inference)

Downsides of CNNs

- Require huge amounts of data and computation for training
- While the individual components are simple, the global behaviour is poorly understood ⇒ lack of guarantees
- Training and fine-tuning is fastidious "Graduate-student descent"
- No free lunch: Trained CNNs are very task/data specific
- Can behave erratically lack of robustness, subject to adversarial attacks
- How the new trend benefits from the techniques of "traditional" IP
 - Deeper understanding of modules
 Suggestion of leaner and more robust architectures

