

EE-559 – Deep learning

7.3. Denoising autoencoders

François Fleuret

<https://fleuret.org/ee559/>

Jan 1, 2020



Beside dimension reduction, autoencoders can capture dependencies between signal components to restore degraded signals.

In that case, we can ignore the encoder/decoder structure, and such a model

$$\phi : \mathcal{X} \rightarrow \mathcal{X}.$$

is referred to as a **denoising** autoencoder.

The goal is not anymore to optimize ϕ so that

$$\phi(X) \simeq X$$

but, given a perturbation \tilde{X} of the signal X , to restore the signal, hence

$$\phi(\tilde{X}) \simeq X.$$

We can illustrate this notion in $2d$ with an additive Gaussian noise, and the quadratic loss, hence

$$\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N \|x_n - \phi(x_n + \epsilon_n; w)\|^2,$$

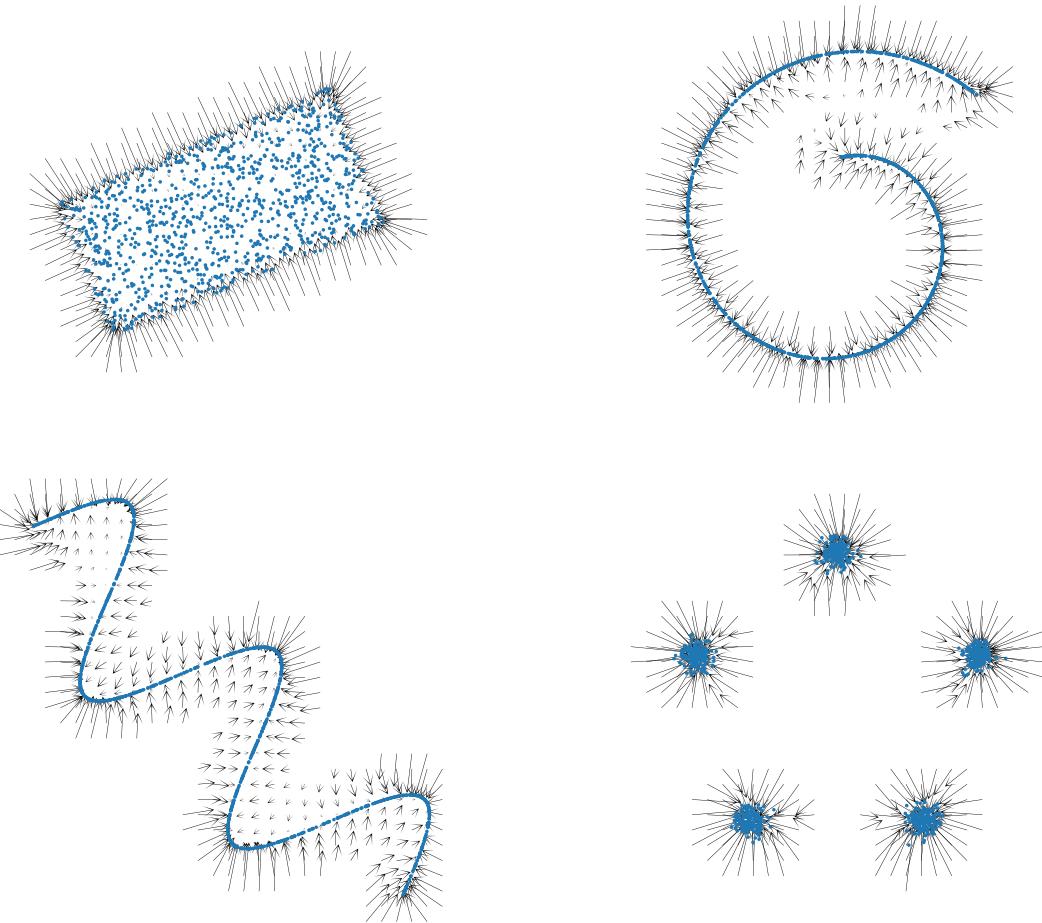
where x_n are the data samples, and ϵ_n are Gaussian random noise vectors.

```
model = nn.Sequential(
    nn.Linear(2, 100),
    nn.ReLU(),
    nn.Linear(100, 2)
)

batch_size, nb_epochs = 100, 1000
optimizer = torch.optim.Adam(model.parameters(), lr = 1e-3)
criterion = nn.MSELoss()

for e in range(nb_epochs):
    for input in data.split(batch_size):
        noise = input.new(input.size()).normal_(0, 0.1)
        output = model(input + noise)
        loss = criterion(output, input)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```



We can do the same on MNIST, for which we keep our deep autoencoder, and ignore its encoder/decoder structure.

```

corrupted_input = corruptor.corrupt(input)

output = model(corrupted_input)
loss = 0.5 * (output - input).pow(2).sum() / input.size(0)

optimizer.zero_grad()
loss.backward()
optimizer.step()

```

We consider three types of corruptions, that go beyond additive noise:

Original



Pixel erasure



Blurring



Block masking



Original

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ($p = 0.9$)

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

7 2 1 0 4 1 4 3 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Original

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Corrupted ($\sigma = 4$)

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Reconstructed

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Original

7	2	/	0	4	/	4	9	5	9	0	6
9	0	1	5	9	7	8	4	9	6	6	5
4	0	7	4	0	1	3	1	3	4	7	2

Corrupted (10×10)

7	2	/	0	4	/	8	5	9	7	0	6
7	0	:	5	9	7	5	4	7	5	6	5
4	0	7	4	0	1	3	1	3	4	7	2

Reconstructed

7	2	/	0	4	/	4	9	5	9	0	6
9	0	1	5	9	7	8	4	9	6	6	5
4	0	7	4	0	1	3	1	3	4	7	2

A key weakness of this type of denoising is that the posterior

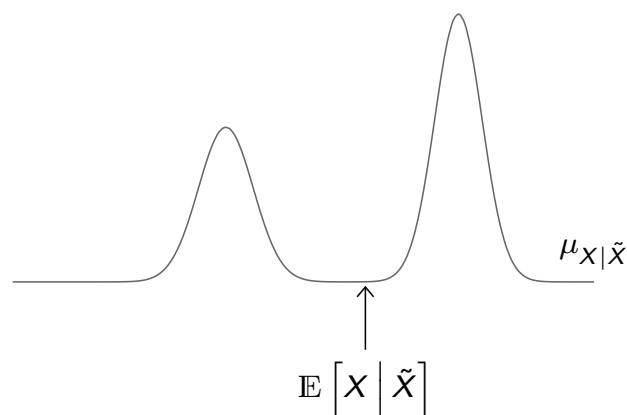
$$\mu_{X|\tilde{X}}$$

may be non-deterministic, possibly multi-modal.

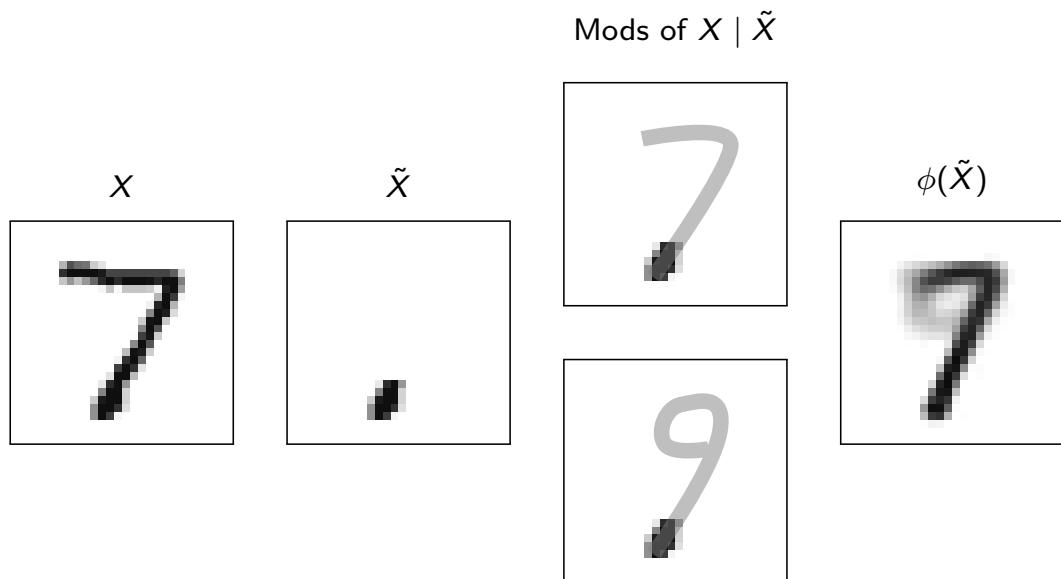
If we train the autoencoder with the quadratic loss, the best reconstruction is

$$\phi(\tilde{X}) = \mathbb{E} [X | \tilde{X}],$$

which may be very unlikely under $\mu_{X|\tilde{X}}$.



This is what we observe with very corrupted MNIST digits.



This can be mitigated by using in place of loss a second network that assesses if the output is realistic.

Such methods are called **adversarial** since the second network aims at spotting the mistakes of the first, and the first aims at fooling the second.

It can be combined with a stochastic denoiser that samples an X according to $X \mid \tilde{X}$ instead of computing a deterministic reconstruction.

We will come back to that.

Noise2Noise

Denoising can be achieved without clean samples, if the noise is additive and unbiased. Consider ϵ and δ two unbiased and independent noises. We have

$$\begin{aligned} & \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - (X + \delta)\|^2 \right] \\ &= \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - X\|^2 \right] - 2\mathbb{E} \left[\langle \delta, \phi(X + \epsilon; \theta) - X \rangle \right] + \mathbb{E} \left[\|\delta\|^2 \right] \\ &= \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - X\|^2 \right] - 2 \left\langle \mathbb{E}[\delta], \mathbb{E}[\phi(X + \epsilon; \theta) - X] \right\rangle + \mathbb{E} \left[\|\delta\|^2 \right] \\ &= \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - X\|^2 \right] + \mathbb{E} \left[\|\delta\|^2 \right]. \end{aligned}$$

Hence

$$\operatorname{argmin}_{\theta} \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - (X + \delta)\|^2 \right] = \operatorname{argmin}_{\theta} \mathbb{E} \left[\|\phi(X + \epsilon; \theta) - X\|^2 \right].$$

Using L_1 instead of L_2 estimates the median instead of the mean, and similarly is stable to noise that keeps the median unchanged.

Lehtinen et al. (2018)'s Noise2Noise approach uses this for image restoration, as many existing image generative processes induce a unbiased noise.

In many image restoration tasks, the expectation of the corrupted input data is the clean target that we seek to restore. Low-light photography is an example: a long, noise-free exposure is the average of short, independent, noisy exposures.

Physically accurate renderings of virtual environments are most often generated through a process known as Monte Carlo path tracing. /.../ The Monte Carlo integrator is constructed such that the intensity of each pixel is the expectation of the random path sampling process, i.e., the sampling noise is zero-mean.

(Lehtinen et al., 2018)

NAME	N_{out}	FUNCTION
INPUT	n	
ENC_CONV0	48	Convolution 3×3
ENC_CONV1	48	Convolution 3×3
POOL1	48	Maxpool 2×2
ENC_CONV2	48	Convolution 3×3
POOL2	48	Maxpool 2×2
ENC_CONV3	48	Convolution 3×3
POOL3	48	Maxpool 2×2
ENC_CONV4	48	Convolution 3×3
POOL4	48	Maxpool 2×2
ENC_CONV5	48	Convolution 3×3
POOL5	48	Maxpool 2×2
ENC_CONV6	48	Convolution 3×3
UPSAMPLE5	48	Upsample 2×2
CONCAT5	96	Concatenate output of POOL4
DEC_CONV5A	96	Convolution 3×3
DEC_CONV5B	96	Convolution 3×3
UPSAMPLE4	96	Upsample 2×2
CONCAT4	144	Concatenate output of POOL3
DEC_CONV4A	96	Convolution 3×3
DEC_CONV4B	96	Convolution 3×3
UPSAMPLE3	96	Upsample 2×2
CONCAT3	144	Concatenate output of POOL2
DEC_CONV3A	96	Convolution 3×3
DEC_CONV3B	96	Convolution 3×3
UPSAMPLE2	96	Upsample 2×2
CONCAT2	144	Concatenate output of POOL1
DEC_CONV2A	96	Convolution 3×3
DEC_CONV2B	96	Convolution 3×3
UPSAMPLE1	96	Upsample 2×2
CONCAT1	$96+n$	Concatenate INPUT
DEC_CONV1A	64	Convolution 3×3
DEC_CONV1B	32	Convolution 3×3
DEV_CONV1C	m	Convolution 3×3 , linear act.

(Lehtinen et al., 2018)

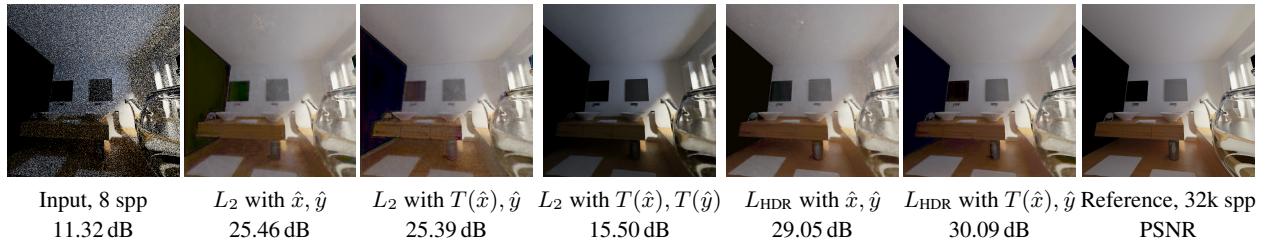


Figure 6. Comparison of various loss functions for training a Monte Carlo denoiser with noisy target images rendered at 8 samples per pixel (spp). In this high-dynamic range setting, our custom relative loss L_{HDR} is clearly superior to L_2 . Applying a non-linear tone map to the inputs is beneficial, while applying it to the target images skews the distribution of noise and leads to wrong, visibly too dark results.



Figure 7. Denoising a Monte Carlo rendered image. (a) Image rendered with 64 samples per pixel. (b) Denoised 64 spp input, trained using 64 spp targets. (c) Same as previous, but trained on clean targets. (d) Reference image rendered with 131 072 samples per pixel. PSNR values refer to the images shown here, see text for averages over the entire validation set.

(Lehtinen et al., 2018)

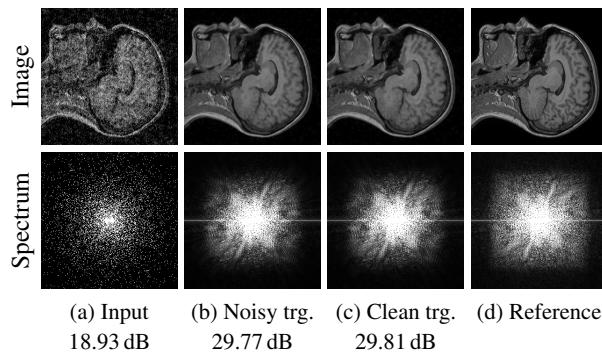


Figure 9. MRI reconstruction example. (a) Input image with only 10% of spectrum samples retained and scaled by $1/p$. (b) Reconstruction by a network trained with noisy target images similar to the input image. (c) Same as previous, but training done with clean target images similar to the reference image. (d) Original, uncorrupted image. PSNR values refer to the images shown here, see text for averages over the entire validation set.

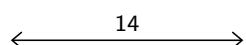
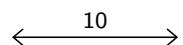
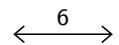
(Lehtinen et al., 2018)

Super-resolution

A special case of denoising is to increase an image resolution. We use an encoder/decoder whose encoder's input is smaller than the decoder's output.

Encoder

Tensor sizes / operations

$1 \times 14 \times 14$	
<code>nn.Conv2d(1, 32, kernel_size=5, stride=1)</code>	 $\xrightarrow{14}$ $\times 10$
$32 \times 10 \times 10$	
<code>nn.Conv2d(32, 32, kernel_size=5, stride=1)</code>	 $\xrightarrow{10}$ $\times 6$
$32 \times 6 \times 6$	
<code>nn.Conv2d(32, 32, kernel_size=4, stride=1)</code>	 $\xrightarrow{6}$ $\times 3$
$32 \times 3 \times 3$	
<code>nn.Conv2d(32, 32, kernel_size=3, stride=1)</code>	 $\xrightarrow{3}$ $\times 1$
$32 \times 1 \times 1$	

```

MNISTUpScaler(
    (encoder): Sequential(
        (0): Conv2d(1, 32, kernel_size=(5, 5), stride=(1, 1))
        (1): ReLU(inplace=True)
        (2): Conv2d(32, 32, kernel_size=(5, 5), stride=(1, 1))
        (3): ReLU(inplace=True)
        (4): Conv2d(32, 32, kernel_size=(4, 4), stride=(1, 1))
        (5): ReLU(inplace=True)
        (6): Conv2d(32, 32, kernel_size=(3, 3), stride=(1, 1))
    )
    (decoder): Sequential(
        (0): ConvTranspose2d(32, 32, kernel_size=(4, 4), stride=(1, 1))
        (1): ReLU(inplace=True)
        (2): ConvTranspose2d(32, 32, kernel_size=(3, 3), stride=(2, 2))
        (3): ReLU(inplace=True)
        (4): ConvTranspose2d(32, 32, kernel_size=(4, 4), stride=(2, 2))
        (5): ReLU(inplace=True)
        (6): ConvTranspose2d(32, 32, kernel_size=(5, 5), stride=(1, 1))
        (7): ReLU(inplace=True)
        (8): ConvTranspose2d(32, 1, kernel_size=(5, 5), stride=(1, 1))
    )
)

```

```

for original, _ in iter(train_loader):
    original = original.to(device)
    input = F.avg_pool2d(original, kernel_size = 2)
    output = model(input)
    loss = (output - original).pow(2).sum() / output.size(0)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

```

Original

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Input

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Bilinear interpolation

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Original

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Input

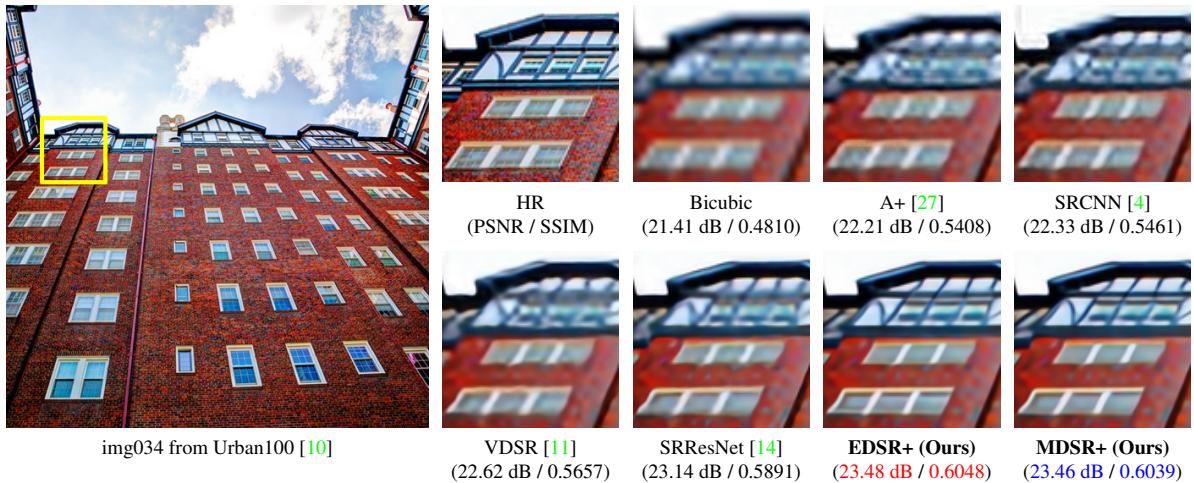
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Autoencoder output

7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 8 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

Lim et al. (2017) use two different resnets.

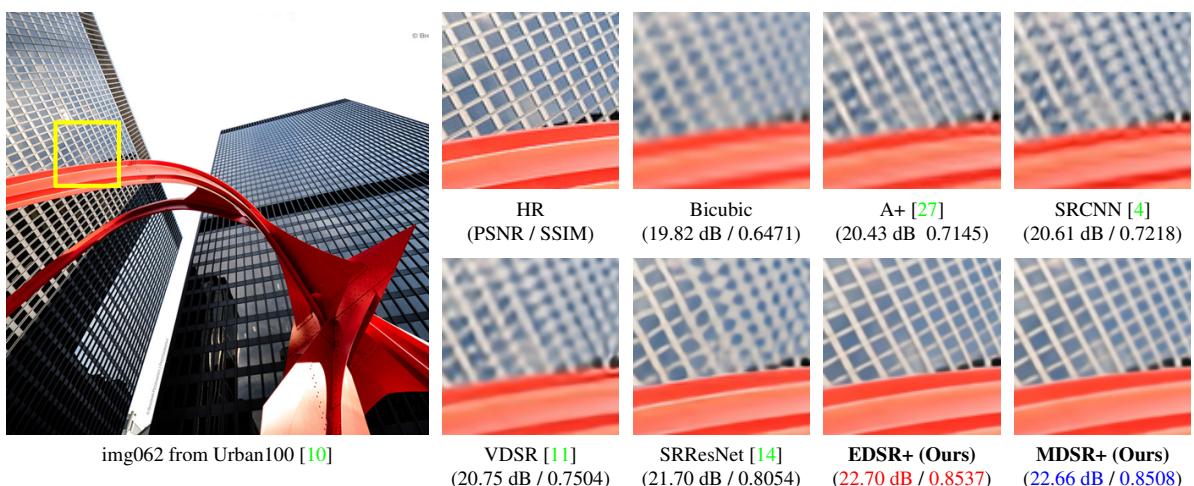
	EDSR	MDSR
Nb blocks	32	80
Channels	256	64
Nb parameters	43M	8M



(Lim et al., 2017)

Lim et al. (2017) use two different resnets.

	EDSR	MDSR
Nb blocks	32	80
Channels	256	64
Nb parameters	43M	8M



(Lim et al., 2017)

Lim et al. (2017) use two different resnets.

	EDSR	MDSR
Nb blocks	32	80
Channels	256	64
Nb parameters	43M	8M



(Lim et al., 2017)

Autoencoders as self-training

Vincent et al. (2010) interpret training the autoencoder as maximizing the mutual information between the input and the latent states.

Let X be a sample, $Z = f(X; \theta)$ its latent representation, and $q(x, z)$ the distribution of (X, Z) .

We have

$$\operatorname{argmax}_{\theta} I(X; Z) = \operatorname{argmax}_{\theta} \mathbb{E}_{q(X, Z)} [\log q(X | Z)].$$

However, there is no expression of $q(X | Z)$ in any reasonable setup.

For any distribution p we have

$$\mathbb{E}_{q(X, Z)} [\log q(X | Z)] \geq \mathbb{E}_{q(X, Z)} [\log p(X | Z)].$$

So we can in particular try to find a “good p ”, so that the left term is a good approximation of the right one.

If we consider the following model for p

$$p(\cdot | Z = z) = \mathcal{N}(g(z; \eta), \sigma)$$

where g is deterministic and σ fixed, we get

$$\mathbb{E}_{q(X, Z)} [\log p(X | Z)] = -\frac{1}{2\sigma^2} \mathbb{E}_{q(X, Z)} [\|X - g(f(X); \eta)\|^2] + k.$$

If optimizing η makes the bound tight, the final loss is the reconstruction error

$$\operatorname{argmax}_{\theta} \mathbb{I}(X; Z) \simeq \operatorname{argmin}_{\theta} \left(\min_{\eta} \frac{1}{N} \sum_{n=1}^N \|x_n - g(f(x_n); \eta)\|^2 \right).$$

This abstract view of the encoder as “maximizing information” justifies its use to build generic encoding layers.

In the perspective of building a good feature representation, just retaining information is not enough, otherwise the identity would be a good choice.

In their work, Vincent et al. propose to degrade the signal with noise before feeding it to the encoder, but to keep the MSE to the original signal.

This makes the encoder retain information about structures beyond local noise.

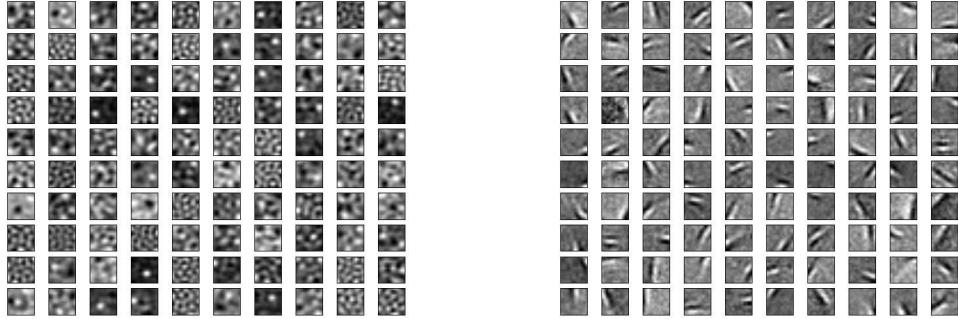


Figure 6: Weight decay vs. Gaussian noise. We show typical filters learnt from natural image patches in the over-complete case (200 hidden units). *Left:* regular autoencoder with weight decay. We tried a wide range of weight-decay values and learning rates: filters never appeared to capture a more interesting structure than what is shown here. Note that some local blob detectors are recovered compared to using no weight decay at all (Figure 5 right). *Right:* a denoising autoencoder with additive Gaussian noise ($\sigma = 0.5$) learns Gabor-like local oriented edge detectors. Clearly the filters learnt are qualitatively very different in the two cases.

(Vincent et al., 2010)

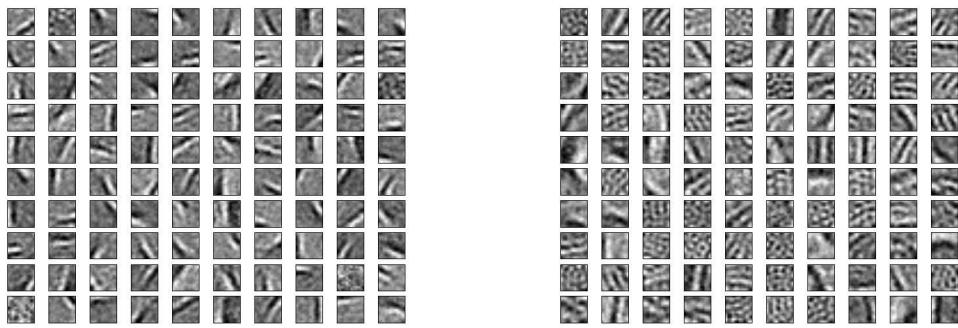


Figure 7: Filters obtained on natural image patches by denoising autoencoders using other noise types. *Left:* with 10% salt-and-pepper noise, we obtain oriented Gabor-like filters. They appear slightly less localized than when using Gaussian noise (contrast with Figure 6 right). *Right:* with 55% zero-masking noise we obtain filters that look like oriented gratings. For the three considered noise types, denoising training appears to learn filters that capture meaningful natural image statistics structure.

(Vincent et al., 2010)

Vincent et al. build deep MLPs whose layers are initialized successively as encoders trained within a noisy autoencoder.

A final classifying layer is added and the full structure can be fine-tuned.

This approach, and others in the same spirit (Hinton et al., 2006), were seen as strategies to complement gradient-descent for building deep nets.

Data Set	SVM _{rbf}	DBN-1	SAE-3	DBN-3	SDAE-3 (v)
MNIST	1.40 ± 0.23	1.21 ± 0.21	1.40 ± 0.23	1.24 ± 0.22	1.28 ± 0.22 (25%)
basic	3.03 ± 0.15	3.94 ± 0.17	3.46 ± 0.16	3.11 ± 0.15	2.84 ± 0.15 (10%)
rot	11.11 ± 0.28	14.69 ± 0.31	10.30 ± 0.27	10.30 ± 0.27	9.53 ± 0.26 (25%)
bg-rand	14.58 ± 0.31	9.80 ± 0.26	11.28 ± 0.28	6.73 ± 0.22	10.30 ± 0.27 (40%)
bg-img	22.61 ± 0.37	16.15 ± 0.32	23.00 ± 0.37	16.31 ± 0.32	16.68 ± 0.33 (25%)
bg-img-rot	55.18 ± 0.44	52.21 ± 0.44	51.93 ± 0.44	47.39 ± 0.44	43.76 ± 0.43 (25%)
rect	2.15 ± 0.13	4.71 ± 0.19	2.41 ± 0.13	2.60 ± 0.14	1.99 ± 0.12 (10%)
rect-img	24.04 ± 0.37	23.69 ± 0.37	24.05 ± 0.37	22.50 ± 0.37	21.59 ± 0.36 (25%)
convex	19.13 ± 0.34	19.92 ± 0.35	18.41 ± 0.34	18.63 ± 0.34	19.06 ± 0.34 (10%)
tzanetakis	14.41 ± 2.18	18.07 ± 1.31	16.15 ± 1.95	18.38 ± 1.64	16.02 ± 1.04 (0.05)

(Vincent et al., 2010)

References

- G. E. Hinton, S. Osindero, and Y.-W. Teh. **A fast learning algorithm for deep belief nets.** *Neural Computation*, 18(7):1527–1554, July 2006.
- J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. **Noise2noise: Learning image restoration without clean data.** *CoRR*, abs/1803.04189, 2018.
- B. Lim, S. Son, H. Kim, S. Nah, and K. Lee. **Enhanced deep residual networks for single image super-resolution.** *CoRR*, abs/1707.02921v1, 2017.
- P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. **Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion.** *Journal of Machine Learning Research (JMLR)*, 11:3371–3408, 2010.