

HW1

Gabriel ROMON

1 DCT denoiser

1. The conditional density of X given $Y = y$ is given by $f_{X|Y=y}(x) = \frac{f_{Y|X=x}(y)f_X(x)}{f_Y(y)}$.
Since $Y = X + B$, $Y|X = x \sim \mathcal{N}(x, \sigma^2)$ hence $f_{X|Y=y}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(y - x)^2)$. Thus

$$f_{X|Y=y}(x) \propto \exp(-\frac{1}{2\sigma^2}(y - x)^2 - |x|)$$

Computing the maximum of $f_{X|Y=y}(x)$ is consequently equivalent to computing

$$\arg \max_{x \in \mathbb{R}} \exp \left(-\frac{1}{2\sigma^2}(y - x)^2 - |x| \right) = \arg \min_{x \in \mathbb{R}} \frac{1}{2\sigma^2}(y - x)^2 + |x|$$

Let $\varphi : x \mapsto \frac{1}{2\sigma^2}(y - x)^2 + |x|$. Studying the derivative of φ on \mathbb{R}^+ and \mathbb{R}^- shows that if $y \geq \sigma^2$, φ is increasing over $[y - \sigma^2, \infty)$ and decreasing on the complement. If $y \leq -\sigma^2$, φ is decreasing over $(-\infty, y + \sigma^2]$ and increasing on the complement. Therefore, if $y \geq \sigma^2$, the minimum is attained at $y - \sigma^2$, if $y \in (-\sigma^2, \sigma^2)$, the minimum is attained at 0 and if $y \leq -\sigma^2$, the minimum is attained at $y + \sigma^2$. This rewrites more compactly as

$$x^* = \left(1 - \frac{\sigma^2}{|y|} \right)^+ y$$

which corresponds to *soft-thresholding* of y with threshold σ^2 .

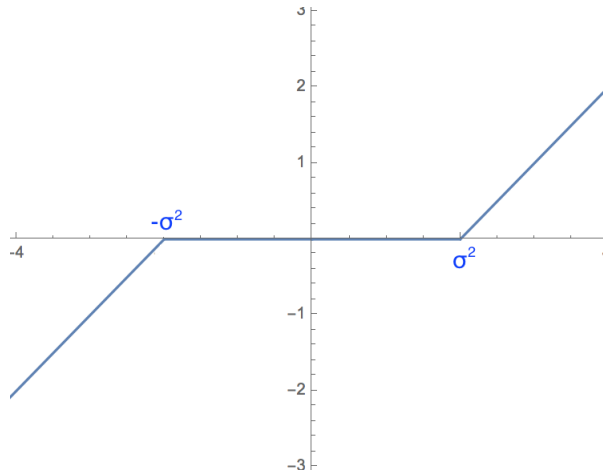


Figure 1: Soft-thresholding function

2. The code in `DCT_denoiser` performs denoising on an image corrupted by Gaussian noise with variance σ^2 .

First, a tensor \mathbf{D} of shape $N \times N \times N^2$ is built such that the $\mathbf{D}[:, :, \mathbf{k}]$ are the inverse discrete cosine transforms of the N^2 elementary matrices of $\mathbb{R}^{N \times N}$. More precisely there is some bijection $\varphi : \llbracket 1, N \rrbracket^2 \rightarrow \llbracket 1, N^2 \rrbracket$ such that $\mathbf{D}[:, :, \varphi(i, j)]$ is the 2D inverse DCT of the elementary matrix E_{ij} . Consequently, $\mathbf{D}[m, n, \varphi(i, j)] = \alpha_i \alpha_j \cos\left(\frac{\pi i}{N}\left(m + \frac{1}{2}\right)\right) \cos\left(\frac{\pi j}{N}\left(n + \frac{1}{2}\right)\right)$ where the α_i are the usual constants defining the type 2 orthonormal DCT.

For some fixed (i, j) , $D[:, :, \varphi(i, j)]$ is a kernel/filter. For each $N \times N$ patch X of the noisy image, `conv2D` computes the correlation between the filter $D[:, :, \varphi(i, j)]$ and X :

$$\sum_{m=0}^{N-1} \sum_{n=0}^{N-1} D[m, n, \varphi(i, j)] X_{mn} = \alpha_i \alpha_j \sum_{m=0}^{N-1} \sum_{n=0}^{N-1} X_{mn} \cos\left(\frac{\pi i}{N} \left(m + \frac{1}{2}\right)\right) \cos\left(\frac{\pi j}{N} \left(n + \frac{1}{2}\right)\right)$$

which is precisely coefficient (i, j) in the 2D DCT of the patch X .

This is done for every (i, j) and the resulting coefficients undergo *hard-thresholding*: all coefficients with absolute value less than 3σ are set to 0. This is different from the soft-thresholding used in 1. Next, the algorithm returns to the spatial domain by convolving the resulting tensor with the same filters (note that this is an actual convolution so filters have to be flipped beforehand, hence the command `np.flipplr(np.flipud(D))`).

3. In order to learn the threshold parameter s , one has to define a loss function and get some training data. Regarding the loss, a possible choice is the squared euclidean norm on the flattened images $\ell(y, \hat{y}) = \sum_i (y_i - \hat{y}_i)^2$. However, the relation between ℓ and s is hard to write down explicitly (even as a composition of several simple functions) thus making gradient descent difficult.
4. Computing the inverse DCT of the elementary matrices has total cost $O(N^2 \times N \log N) = O(N^3 \log N)$. For an $m \times n$ input, the correlation operation has cost $O(N^2 \times (m - N)(n - N) \times N^2) = O(mnN^4)$ since for each filter, there are $\sim (m - N)(n - N)$ patches and $O(N^2)$ operations to compute on each. The second convolution has the same cost $O(mnN^4)$, so the number of operations per pixel is $O(N^4)$.

2 DnCNN

5. The architecture used in the practical session is identical to the one presented in the seminal paper on DnCNN [1]: the network has depth 17, the first layer is Conv+ReLU (64 filters of size $3 \times 3 \times 1$), each layer from №2 to №16 is made of Conv+BatchNorm+ReLU (64 filters of size $3 \times 3 \times 64$) and the final layer is convolutional (1 filter of size $3 \times 3 \times 64$).
Layer №1 has $3 \times 3 \times 1 \times 64 + 64$ parameters, while layers from №2 to №16 each have $3 \times 3 \times 64 \times 64 + 64$ parameters and the last has $3 \times 3 \times 64 \times 1 + 1$, for a grand total of ~ 555.000 parameters (weights and biases).
6. For an $m \times n$ input, convolutions in layer №1 take $\sim 64 \times 3^2 \times mn$ operations, convolutions in each of the layer №2 to №16 take $\sim 64^2 \times 3^2 \times mn$ and for the last layer $64 \times 3^2 \times mn$. For the network used in the practical session, the number of operations per pixel is ~ 554.000 .

Let D denote the depth of the network ($D = 17$ in the implementation). Since the network stacks convolutional layers of 3×3 filters, the second layer has a 3×3 view of the first one, hence a 5×5 view of the input, and so on with the third layer which has a 7×7 view of the input. The final layer has a receptive field of size $(2D + 1) \times (2D + 1)$, which is exactly what [1] states. If one wants to compare DCN denoising with DnCNN, one should thus compare the patch size N with $2D + 1$ (and not with 3, the size of the filters).

Throwing away constants, the total number of operations is thus $O(Dmn)$, which yields $O(D)$ operations per pixel. Given what has been said above, N can be compared with D ($2D + 1$ actually but this is harmless) and the complexity has changed from $O(N^4)$ to $O(D)$, a significant improvement, but one should not forget that DnCNN needs to be trained on a large dataset beforehand.

7. Figure 1 shows the performance of both methods on a picture corrupted by Gaussian noise (with standard deviation $\sigma = 25$, the same used for training the network). DCT denoising performs correctly, but the resulting image shows multiple significant local artifacts (as if its resolution had decreased). DnCNN returns a picture that is more contrasted and more detailed compared to the other method (zoom in to see the difference).



DCT Denoising



DnCNN Denoising

Figure 2: Denoised images obtained with each method

8.

9. Figure 2 shows the performance of both methods on a picture corrupted by Gaussian noise with standard deviation $\sigma = 5$. The image returned by DCT denoising (after adjusting the threshold to $s = 15$) is very similar to the original, with a minor loss of details. DnCNN returns an image without visible noise, but it is much dimmer (low contrast), grass has disappeared (meaning that high frequencies are lost) and there are significant artifacts.



DCT Denoising



DnCNN Denoising

Figure 3: Denoised images obtained with each methods for $\sigma = 5$

3 DRCN and DCSCN

10. According to [2] DRCN is made up of three sub-networks. First, an embedding net with 2 stacked Conv+ReLU takes as input an interpolated version of the image and computes feature maps. Then, those are fed to an inference net that stacks 9 recursive layers, each made of the **same** convolution layer followed by ReLU. This parameter sharing between layers helps to reduce the complexity of the model. Finally, the high-resolution image is built by a reconstruction net: it is connected directly to the original input (skip-connection is useful when the input and output are highly correlated) and to each of the recursive layers (this helps mitigate vanishing and exploding gradients). According to [3] DCSCN takes as input the original image (not upsampled or interpolated) and begins with a feature extraction network that stacks 7 Conv+PReLU with decreasing depth. The output of each of these layers are concatenated in a single matrix. This matrix is then fed to a reconstruction network that learns the high resolution image minus the bicubic interpolation of the input (residual learning similar to DnCNN).

11. Regarding DRCN, the convolutional layers in the practical session use $3 \times 3 \times 1$ and $3 \times 3 \times 96$ filters. The number of operations per pixel is

$$9 \cdot 96 + 9 \cdot 96^2 + 9^2 \cdot 96^2 + 9(9 \cdot 96^2 + 9 \cdot 97) + 9 \cdot 4 \sim 1.584.000$$

DCSCN has components of varying sizes. The number of operations per pixel is

$$9 \cdot (32 + 32 \cdot 26 + 26 \cdot 22 + 22 \cdot 18 + 18 \cdot 14 + 14 \cdot 11 + 11 \cdot 8) + 131 \cdot 24 + 131 \cdot 8 + 9 \cdot 8 \cdot 8 + 9 \cdot 32 \cdot 128 + 9 \cdot 4 \cdot 32 \sim 64.000$$

12. Input images have been rescaled by a factor $\frac{1}{2}$ using two methods. The brutal method simply drops every other row and every other column. The PIL method uses bicubic interpolation after downsampling. As seen in Figure 3, the first technique yields an image that looks more pixelated.

Super-resolution results are in Figure 4. DRCN is clearly affected by the downsampling technique since the details in the high-resolution image remain very blurry. This may be explained by the training process because DRCN takes interpolated images as input. On the contrary the outputs of DCSCN look similar.



Brutal dezoom



PIL dezoom

Figure 4: Zoomed out test images



Brutal and DRCN



PIL and DRCN



Brutal and DCSCN



PIL and DCSCN

Figure 5: High-resolution images obtained with both methods

13. For the PIL image, intermediate images produced by the reconstruction net are shown in Figure 5. They are taken from the output of the reconstruction net, just before the averaging step. I do not see any significant difference or visible trend in these images.

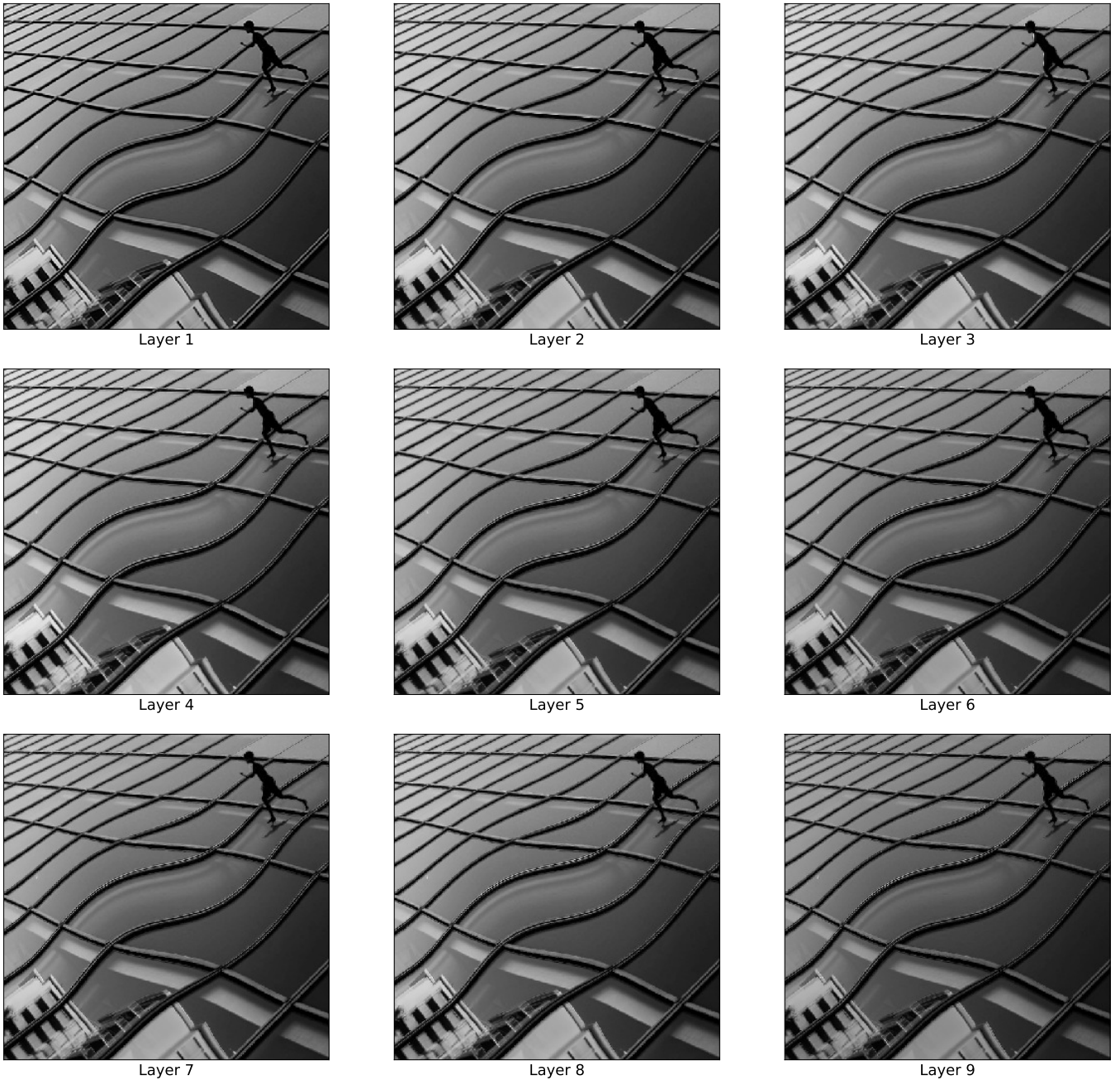


Figure 6: Intermediate images produced by the reconstruction net (before averaging)

4 Some training

14. The DCSCN architecture has been modified. The network takes as input the original high-resolution image, convolves it with a kernel, downsamples it and feeds it to the usual DCSCN architecture. The goal is to learn the kernel by gradient descent (the DCSCN comes pretrained).

References

- [1] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017.
- [2] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016.
- [3] Jin Yamanaka, Shigesumi Kuwashima, and Takio Kurita. Fast and accurate image super resolution by deep cnn with skip connection and network in network. In *Neural Information Processing*, pages 217–225. Springer, 2017.