

# Can a Single Image Denoising Neural Network Handle All Levels of Gaussian Noise?

Yi-Qing Wang, and Jean-Michel Morel

**Abstract**—A recently introduced set of deep neural networks designed for the image denoising task achieves state-of-the-art performance. However, they are specialized networks in that each of them can handle just one noise level fixed in their respective training process. In this note, by investigating the distribution invariance of the natural image patches with respect to linear transforms, we show how to make a single existing deep neural network work well across all levels of Gaussian noise, thereby allowing to significantly reduce the training time for a general-purpose neural network powered denoising algorithm.

**Index Terms**—deep neural network, distribution invariance, image denoising, natural patch space

## I. INTRODUCTION

RECENTLY, a set of deep neural networks, or multi-layer perceptrons (MLPs) designed for the image denoising task [1] has been shown to outperform BM3D [2], widely accepted as the state-of-the-art. Although these enormous deep networks only work at the noise levels which they were trained for, this turn of events clearly demonstrates the potential of a pure learning strategy. A philosophical difference sets the two patch-based methods apart: BM3D, a major spin-off of the original non-local means [3], seeks information exclusively inside the noisy image while the neural network derives all its power by looking at noisy and clean patch pairs gathered from other images. Other algorithms exist [4]–[8] which fall between these two ends of the spectrum.

It is well known [9], [10] that neural networks form a class of universal approximators. Recent studies [11] further suggest that multi-layer neural networks tend to be more efficient in signal representation than their traditional one-hidden-layer counterpart. Whatever their architectures, neural networks in a regression framework seek to approximate the conditional expectation under some input distribution. In our setting, let  $x$ ,  $\tilde{x}$ , and  $y$  denote the clean, noisy and denoised patch. Note that  $y$  does not necessarily have the same dimension as  $\tilde{x}$ . All the neural networks [1] under this study, for instance, produce a  $y \in \mathbb{R}^{17 \times 17}$  based on a noisy observation  $\tilde{x} \in \mathbb{R}^{39 \times 39}$  which includes not only  $y$ 's corresponding noisy pixels but also its surrounding ones. Also note that with the *true* natural patch distribution beyond reach, a huge number of patches are drawn stochastically [12] from a large natural image dataset

to provide the underlying distribution for computing

$$\begin{aligned} \theta^* &= \operatorname{argmin}_{\theta} \mathbb{E} \|f(\tilde{x}, \theta) - x\|_2^2 \\ &= \operatorname{argmin}_{\theta} \mathbb{E} \|f(\tilde{x}, \theta) - \mathbb{E}[x|\tilde{x}]\|_2^2 \end{aligned} \quad (1)$$

where  $f(\cdot, \theta) : \tilde{x} \mapsto y$  is a neural network parametrized by its connection weights  $\theta$ . Due to this problem's non-convex nature in general, instead of the potentially intractable  $\theta^*$ , a good  $\theta$ , judged on the basis of the resulting network's generalization error, is usually accepted as a solution.

In spite of their impressive performance, the proposed deep networks are impractical. As stated in the paper itself [1]: *our most competitive MLP is tailored to a single level of noise and does not generalize well to other noise levels. This is a serious limitation which we already tried to overcome with an MLP trained on several noise levels. However, the latter does not yet achieve the same performance for  $\sigma = 25$  as the specialized MLP.* Since a standard general-purpose algorithm for Gaussian noise removal ought to be able to handle all levels of noise, this limitation seems to require a series of such networks, one for each noise level, which is impractical and even unrealistic especially in view of their prohibitive training time [13].

In this note, through an analysis of the interplay between the neural networks and the underlying patch distribution they seek to learn, we show how to construct a linear transform that moves natural image patches within the support of their distribution, which is then used to make a single existing deep neural network work well across all levels of Gaussian noise. In the concluding section, using the natural patch distribution invariance argument, we hint that further patch normalization may help scale down these deep neural networks by reducing their domain of definition without compromising their power.

## II. A SINGLE NEURAL NETWORK FOR ALL NOISE LEVELS

To make a single neural network work for all noise levels, first we need to investigate the statistical regularity of the *natural patch space*, or the support of the natural patch distribution, with respect to some linear transforms: we drew  $10^6$  39-by-39 random patches from the *Berkeley Segmentation Dataset* (BSD500) rendered to grayscale with *Matlab*'s `rgb2gray` function. The patches were then normalized using the formula

$$\bar{p} = (p/255 - 0.5) \cdot 5 \quad (2)$$

provided in [1] in order to conform them to these deep neural networks' training patch distribution. For each normalized patch, we computed the mean and standard deviation of its 1521 pixels and then plotted their population distribution along

This work was supported in part by DxO Labs, ERC(AG Twelve Labours) and ONR N00014-97-1-0839.

The authors are with CMLA, Ecole Normale Supérieure de Cachan, Cachan, 94230, France (e-mail: yqwang9@gmail.com; morel@cmla.ens-cachan.fr).

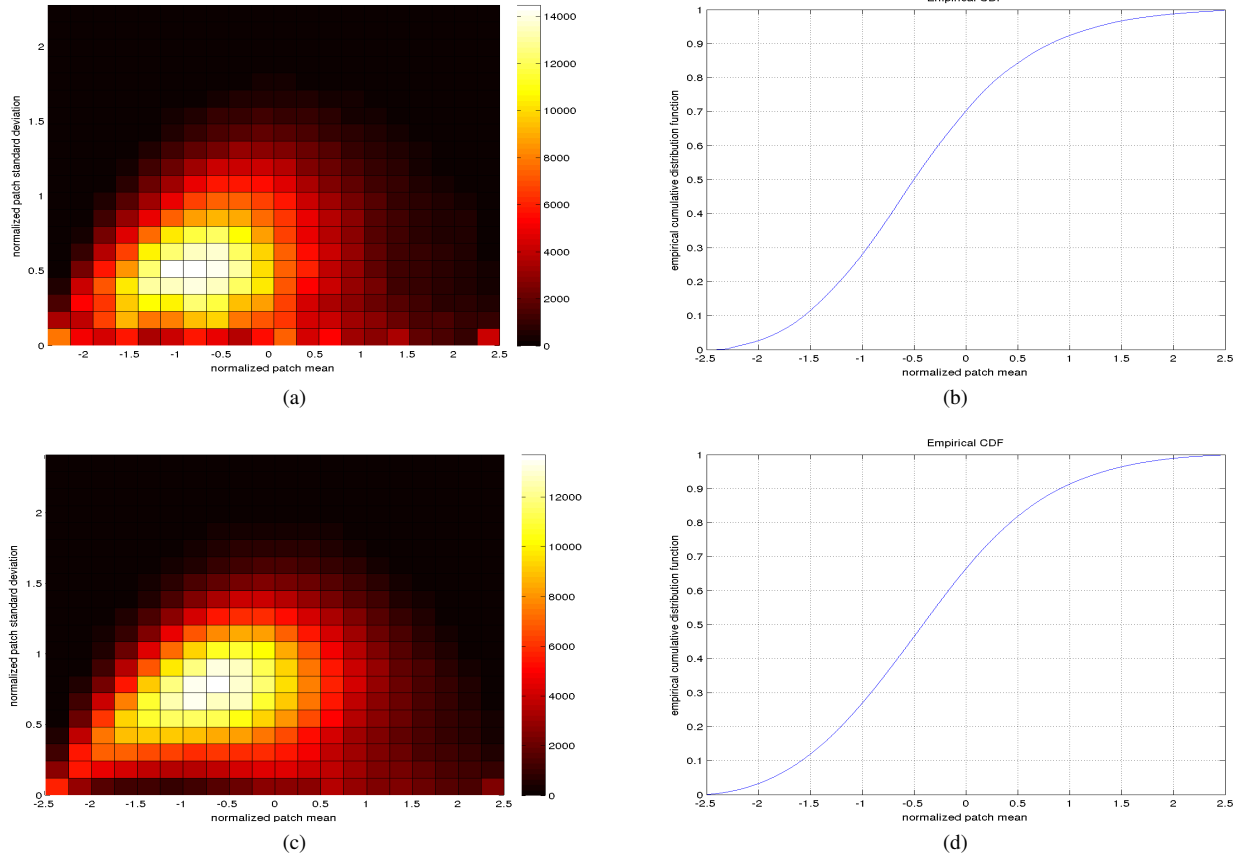


Fig. 1. 2D histogram of random patches from (1a) the grayscale BSD500 and (1c) the grayscale *PASCAL VOC 2012*. Their horizontal (resp. vertical) axis represents the normalized patch's mean (resp. standard variation). (1d) and (1b) plot their respective marginal cumulative distribution function along the patch mean. In both cases, the normalized patches have their mean concentrated around  $-0.5$ . Also note that at the two ends of the horizontal axis, there are two important flat patch clusters because of saturation.

these two dimensions. For comparison, we did the same with the *PASCAL VOC 2012* dataset.

The results (Fig.1) show that the means of the normalized natural patches concentrate around  $-0.5$ . Even without access to the supervised pairs used in [1], we can therefore expect their neural networks trained according to the criterion (1) to do well with patches from this neighborhood because of the sheer number of examples available. Moreover, the patch-wide variance peaking at the patch-wide means around  $-0.5$  strongly indicates a higher tolerance for linear transforms there, that is, it is more probable for a natural patch  $p$  transformed by

$$q = a \cdot (p - s) \text{ with } s \text{ some constant patch and } a \geq 0$$

to remain *natural* if  $q$ 's mean is close to  $-0.5$ .

To verify this conjecture that a linear transform exists which only moves patches within the support of the natural patch distribution, we designed a test that shifts the means of the patches to the same value before denoising them using a neural network, after which the shifted differences were added back individually to obtain their final, denoised versions. It is important to note that thanks to a high ratio between the patch size and  $\sigma$ , estimating the clean patch mean from its noisy version (Line 7 in Test) is very reliable.

---

### Test Patch Mean Normalization Test

---

- 1: **Input:**  $n$  clean patches  $p_j$  of dimension  $39 \times 39$
  - 2: **Output:**  $n$  denoised (resp. clean)  $17 \times 17$  patches  $\hat{p}_j$  (resp.  $p_j$ )
  - 3: **Parameter:** noise variance  $\sigma^2$  and desired patch mean value  $m$
  - 4: **for**  $j = 1$  to  $n$  **do**
  - 5:     retrieve the  $17 \times 17$  clean patch  $p_j$  in the center of  $p_j$
  - 6:     generate the normalized noisy patch
 
$$\tilde{x}_j \leftarrow ((p_j + n_j)/255 - 0.5) \cdot 5$$
 with  $n_j$  having 1521 i.i.d. Gaussian random variables  $\mathcal{N}(0, \sigma^2)$  and  $n_j$  independent of  $n_{j'}$  for  $j \neq j'$
  - 7:     compute the patch mean shift  $s_j = \frac{1}{1521} \sum_{k=1}^{1521} \tilde{x}_{jk} - m$  where  $\tilde{x}_{jk}$  is the  $k$ -th pixel in  $\tilde{x}_j$
  - 8:     shift the network's input  $\forall k, \tilde{x}_{jk} \leftarrow \tilde{x}_{jk} - s_j$
  - 9:     denoise  $y_j = f_\sigma(\tilde{x}_j, \theta)$  where  $f_\sigma(\cdot, \theta)$  is the deep neural network [1] trained with Gaussian noise standard deviation  $\sigma$
  - 10:     shift the network's output by the same amount  $\forall k, y_{jk} \leftarrow y_{jk} + s_j$  in the opposite direction
  - 11:     reverse the normalization  $\hat{p}_j \leftarrow (y_j/5 + 0.5) \cdot 255$
  - 12: **end for**
- 

Running the test on  $10^5$  39-by-39 random patches drawn from the grayscale BSD500, we computed the resultant root mean square error (RMSE). Fig.2 shows that regardless of

the applied noise strength  $\sigma$ , the best empirical performance is attained with the patch mean shift set to  $-0.5$ , thereby fully confirming our supposition. In addition, observe that the optimal patch mean shift incurs surprisingly little RMSE loss, which is less than 0.1 for a noise level as high as 75.

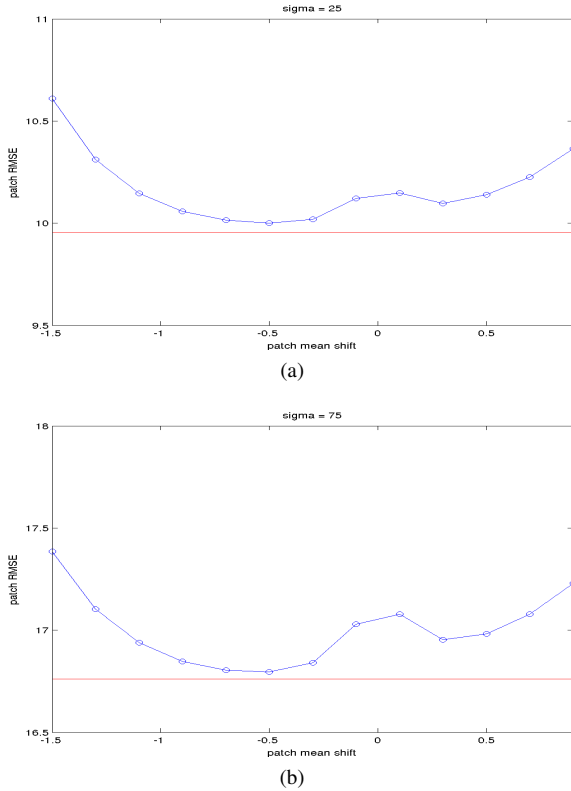


Fig. 2. (2a) (resp. (2b)) plots the RMSEs resulting from the test with  $\sigma = 25$  (resp. 75). The blue curve records the dedicated network's performance with the patch mean values ranging from  $-1.5$  to  $0.9$ . And the red line is the RMSE achieved on the same test data without patch mean shift. Other available neural networks have also been tested with very similar results.

The previous analysis paves the way for a generic network able to handle all levels of noise. Let us use the neural network trained with the noise level  $\sigma^*$  for instance. To restore a noisy patch  $\tilde{p}$  with noise standard deviation at  $\sigma$ , one multiplies  $\tilde{p}$  by  $\sigma^* \sigma^{-1}$ , apply the patch mean shift and let the neural network operate on the normalized patch (see Algorithm). Henceforth  $\sigma^*$  itself becomes a parameter for further optimization, too.

A comparison among various  $\sigma^*$ -indexed generic networks constructed this way (see Fig.3) shows that the one built with  $\sigma^* = 25$  is the best, which is hardly surprising because a neural network can learn the most about the underlying patch space when noise is weak. Moreover, except for the extremely noisy case, one sees no tangible difference in Fig.3b between a dedicated network and the best generic network. Also observe that the higher the  $\sigma^*$ , the worse the resulting RMSEs at low noise levels, an expected phenomenon since a high  $\sigma^* \sigma^{-1}$  exaggerates the patch-wide variation so much that resultant patches no longer remain in the natural patch space.

The implication of this analysis for future research is straightforward: training a network with the same architecture but at an even lower noise level may be rewarding. But it should also be said that too low a  $\sigma^*$  is not likely to work well

### Algorithm Generic Neural Network Denoising

- 1: **Input:** noisy patch  $\tilde{p}$  of dimension  $39 \times 39$  and its noise level  $\sigma$
- 2: **Output:** denoised  $17 \times 17$  patch  $\hat{p}$
- 3: **Parameter:** noise level  $\sigma^*$  of the trained neural network  $f_{\sigma^*}(\cdot, \theta)$  and optimal shift  $m = -0.5$
- 4: Scale the noise  $\tilde{p} \leftarrow \sigma^* \sigma^{-1} \tilde{p}$
- 5: Normalize the patch  $\tilde{x} \leftarrow (\tilde{p}/255 - 0.5) \cdot 5$
- 6: Compute the patch mean shift  $s = \frac{1}{1521} \sum_{k=1}^{1521} \tilde{x}_k - m$  where  $\tilde{x}_k$  is the  $k$ -th pixel in  $\tilde{x}$
- 7: Shift the network's input mean  $\forall k, \tilde{x}_k \leftarrow \tilde{x}_k - s$
- 8: Run the generic neural network denoising  $y = f_{\sigma^*}(\tilde{x}, \theta)$
- 9: Shift the network's output by the same amount  $\forall k, y_k \leftarrow y_k + s$  in the opposite direction
- 10: Reverse the normalization  $\hat{p} \leftarrow (y/5 + 0.5) \cdot 255 \cdot (\sigma^* \sigma^{-1})^{-1}$

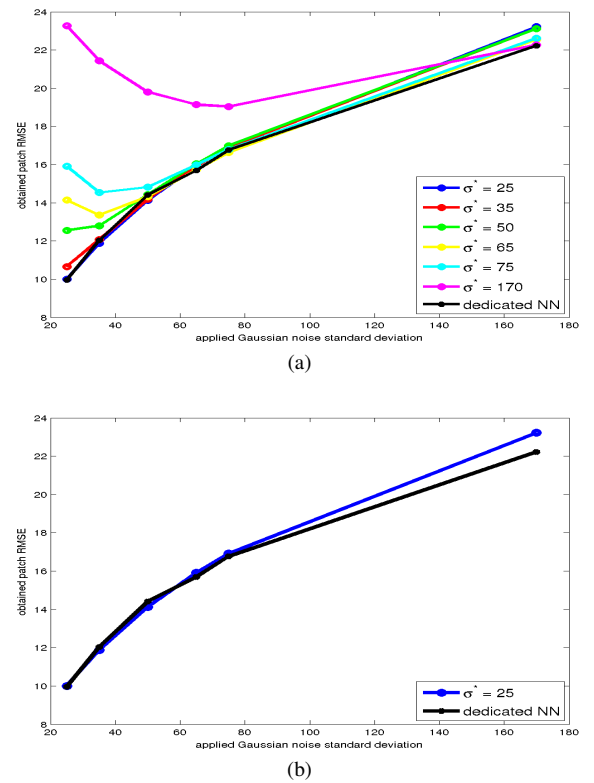


Fig. 3. (3a) Performance discrepancy between dedicated neural networks and  $\sigma^*$ -indexed generic neural networks at handling different noise levels. The horizontal axis marks various test Gaussian noise levels and the vertical axis the achieved RMSEs on  $10^5$  random patches from BSD500. (3b) singles out the best generic network with  $\sigma^* = 25$ .

in a strong noise environment, as already observed in Fig.3. Because the factor  $\sigma^* \sigma^{-1}$  will then flatten all the meaningful patch-wide variations, leaving the network unable to tell one patch from another.

Put differently, what we have shown is that the trained deep network is not very different from its most informative section. Hence, one may want to train a network on mean normalized patches to improve performance, thanks to a denser data distribution. However, in so doing, one implicitly trains on a marginal distribution and thus risks losing information.

To conclude, we tested the Algorithm with  $\sigma^* = 25$  on some

TABLE I  
COMPARISON BETWEEN DEDICATED NEURAL NETWORKS AND OUR  
ALGORITHM (GENERIC  $\sigma^* = 25$ )

$\sigma = 25$	dedicated	generic	$\sigma = 35$	dedicated	generic
computer	<b>7.99</b>	8.10	computer	<b>9.63</b>	9.82
dice	2.77	2.77	dice	<b>3.29</b>	3.45
flower	<b>4.59</b>	4.63	flower	<b>5.53</b>	5.64
girl	<b>3.38</b>	3.41	girl	<b>3.88</b>	4.07
traffic	<b>9.16</b>	9.22	traffic	<b>10.81</b>	10.94
valldemossa	<b>11.85</b>	11.99	valldemossa	<b>14.21</b>	14.28
avg.	<b>6.62</b>	6.68	avg.	<b>7.89</b>	8.03

$\sigma = 50$	dedicated	generic	$\sigma = 65$	dedicated	generic
computer	<b>11.59</b>	11.92	computer	<b>13.16</b>	13.82
dice	<b>4.17</b>	4.50	dice	<b>4.83</b>	5.39
flower	<b>6.85</b>	7.06	flower	<b>7.89</b>	8.20
girl	<b>4.60</b>	4.92	girl	<b>5.28</b>	5.74
traffic	<b>12.83</b>	13.07	traffic	<b>14.32</b>	14.78
valldemossa	<b>16.98</b>	17.06	valldemossa	<b>18.67</b>	18.99
avg.	<b>9.50</b>	9.75	avg.	<b>10.69</b>	11.15

$\sigma = 75$	dedicated	generic	$\sigma = 170$	dedicated	generic
computer	<b>14.10</b>	14.78	computer	<b>20.51</b>	21.81
dice	<b>5.48</b>	6.14	dice	<b>9.23</b>	10.92
flower	<b>8.57</b>	8.96	flower	<b>12.84</b>	13.64
girl	<b>5.66</b>	6.20	girl	<b>8.79</b>	10.54
traffic	<b>15.08</b>	15.56	traffic	<b>20.71</b>	21.72
valldemossa	<b>19.97</b>	20.34	valldemossa	<b>26.42</b>	27.26
avg.	<b>11.47</b>	11.99	avg.	<b>16.41</b>	17.64

real images. The test set (Fig.4) comprises six noiseless images proposed for benchmarking various denoising algorithms [14] and the results are compiled in Tab.I. Recall that even for  $\sigma = 25$ , our generic network is different from the dedicated one, in that ours involves one additional step of patch mean shift. The obtained results are thus consistent with Fig.2a.



Fig. 4. Test images for algorithm comparison (4a) computer (4b) dice (4c) flower (4d) girl (4e) traffic (4f) valldemossa. All images are of dimension  $704 \times 469$  except for valldemossa ( $769 \times 338$ ).

### III. CONCLUSION

In this note, we have shown how to make a single existing neural network work well across all levels of Gaussian noise, thereby allowing to reduce significantly the training time for a general-purpose neural network powered denoising algorithm.

To make deep neural network based algorithm more practical, the other major challenge is to reduce their sizes. This might be achieved through further patch normalization so as to reduce the neural network's input complexity. Rotation and scale invariances of natural image statistics might be used for that purpose.

### REFERENCES

- [1] H. Burger, C. Schuler, and S. Harmeling, "Image denoising: Can plain neural networks compete with BM3D?" in *IEEE Conf. Computer Vision and Pattern Recognition*, 2012, pp. 2392–2399.
- [2] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image restoration by sparse 3D transform-domain collaborative filtering," in *Electronic Imaging*, 2008.
- [3] A. Buades, B. Coll, and J. M. Morel, "A review of image denoising algorithms, with a new one," *Multiscale Modeling & Simulation*, vol. 4, no. 2, pp. 490–530, 2005.
- [4] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans. Image Processing*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [5] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Non-local sparse models for image restoration," in *IEEE Int. Conf. Computer Vision*, 2009, pp. 2272–2279.
- [6] D. Zoran and Y. Weiss, "From learning models of natural image patches to whole image restoration," in *IEEE Int. Conf. Computer Vision*, 2011, pp. 479–486.
- [7] Y. Q. Wang and J. M. Morel, "SURE guided gaussian mixture image denoising," *SIAM Journal on Imaging Sciences*, vol. 6, no. 2, pp. 999–1034, 2013.
- [8] H. Burger, C. Schuler, and S. Harmeling, "Learning how to combine internal and external denoising methods," in *Pattern Recognition*. Springer, 2013, pp. 121–130.
- [9] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [10] A. R. Barron, "Approximation and estimation bounds for artificial neural networks," *Machine Learning*, vol. 14, no. 1, pp. 115–133, 1994.
- [11] Y. Bengio, "Learning deep architectures for ai," *Foundations and Trends® in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [12] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. Int. Conf. Computational Statistics*. Springer, 2010, pp. 177–186.
- [13] H. Burger, "Modelling and learning approaches to image denoising," Ph.D. dissertation, Eberhard Karls Universität Tübingen, Wilhelmstr. 32, 72074 Tübingen, 2013.
- [14] M. Lebrun, "An Analysis and Implementation of the BM3D Image Denoising Method," *Image Processing On Line*, vol. 2012, pp. 175–213, 2012.