

Image Denoising survey

Monday, November 4, 2024 8:59 AM

paper: Image Denoising: The Deep Learning Revolution and Beyond – A Survey Paper by Michael Elad, Bahjat Kawar, Gregory Vaksman
SIAM J. Imaging Science, 2023

① What is image denoising?

$$y = x + v, \quad \text{all vectors in } \mathbb{R}^N \quad (\text{eg. } N = n_x \times n_y)$$

↑ ↑ ↗
observed true image noise

* note we will use vectors not matrices

It's the simpliest inverse problem ($y = A \cdot x + v$) ie. $A = I$

- Assume $v \sim N(0, \sigma^2 I)$
 ↗ and σ^2 known.
- } "AWGN"
 } Additive, white Gaussian Noise

... the simplest noise model.

→ real cameras (CCD pixel arrays) have a tiny amount of AWGN. Mostly "shot noise" due to quantum nature of discrete photons, following a Poisson distribution (partly addressed via Anscombe or other variance stabilizing distribution, model as Gaussian if photon count high)

• Why so simple?

- Fundamental: backbone of more realistic setups
- } indirect
($A \neq I$, and other noise)

- Tractable, understand

- Paper argues its a key ingredient in fancier
- } direct
setups

• Goal

Build a denoiser D , return estimate

$$\hat{x} = D(y, \sigma)$$

Image Denoising survey (p. 2)

Monday, November 4, 2024 9:11 AM

- What criterion?

For now, for good reason, use MSE (equivalently, PSNR)

* Careful: before, "MSE" is $\frac{1}{\text{#pixels}} \sum_{i=1}^{\text{#pixels}} (\hat{x}_i - x_i)^2$
and "mean" referred to this average
actual constant of little importance if it never changes

Now,

$$\text{MSE} = \mathbb{E} \|\hat{x} - x\|_2^2 := \int \|\hat{x} - x\|_2^2 \cdot p(x) dx$$

$$(\hat{x} = D(y, \sigma), y = x + v, \text{ so } \hat{x} \text{ is a function of } x)$$

i.e. "mean" is used in a different sense.

That is, we're assigning a prior distribution to the set of images.

$$p(X = \boxed{\text{house}}) = 0.013, \quad p(X = \boxed{\text{noisy image}}) = 10^{-20}$$

This is the critical ingredient...
without a prior, there's no hope

Setting $\hat{x} = y$ is the best you can do.

So this is Bayesian?
not Frequentist?

Note to future years:
lecture was 1 day before
2024 presidential election

Yes, sort of...

but we don't (anymore) write down a formula for $p(x)$
(these days, we'll learn it from our dataset)

Image Denoising survey (p. 3)

Monday, November 4, 2024 9:21 AM

Estimation 101

$$y = x + v, \quad P_x \text{ is prior on } X$$

$P_{x,y}$ is joint on $X \times Y$

whether Bayesian or not, Bayes' rule is true:

$$P_{X|Y=y}(x) = \frac{P_{x,y}(x,y)}{P_y(y)} = \frac{P_{y|x=x}(y) P_x(x)}{P_y(y)}$$

" x " refers to r.v.

" x " refers to possible value of r.v.

Maximum Likelihood Estimation (MLE) [not "Bayesian"] Very popular

Find x to maximize $p(y|x)$ ("likelihood")

Given x , $y = x + v$ i.e. $y \sim N(x, \sigma^2 I)$

$$\text{so } p(y|x) = \text{const. } \exp(-\frac{1}{2} \frac{1}{\sigma^2} \|y - x\|^2)$$

so

$$\hat{x}_{\text{MLE}} = \underset{x}{\operatorname{argmin}} -\log(p(y|x)) = \underset{x}{\operatorname{argmin}} \frac{1}{2} \frac{1}{\sigma^2} \|y - x\|^2 = y$$

$$\hat{x}_{\text{MLE}} = y \text{ not useful! (if } A \neq I \text{ it can be useful)}$$

Maximum 'a posteriori' Estimation (MAP) [Bayesian]

Find x to maximize $p(x|y)$

$$p(x|y) = \frac{p(y|x) p(x)}{p(y)}$$

unimportant

$$\min_x -\log(p(x|y)) = \min_x -\log(p(y|x)) - \log(p(x))$$

MLE term

$p(x)$

if $p(x) = \text{constant}$ (i.e. uniform, uninformative)

then this has no effect
So MLE=MAP
in this case

exploits a prior

$$\text{assume } p(x) \sim e^{-\rho(x)}$$

data fitting term

regularization term

Image Denoising survey (p. 4)

Monday, November 4, 2024 11:02 AM

Ex Suppose $p(x) = \text{const} \cdot \exp(-\frac{\rho(x)}{2})$ i.e. we prioritize small values of pixels

then MAP:

$$\min_x -\log(p(y|x)) - \log(p(x))$$

$$= \min_x \frac{1}{2\sigma^2} \|x-y\|^2 + \frac{1}{2} \|x\|^2 \quad \text{Tikhonov regularization / Ridge regression}$$

MMSE minimum MSE estimator [Bayesian]

Find \hat{x} to minimize $MSE \mathbb{E} \|\hat{x}-x\|^2$

$$\hat{x}_{\text{MMSE}} = \mathbb{E}[x|y] \quad \text{"closed form" misleading not helpful directly except in simplest cases}$$

(2) "Classical" Techniques

... most (not all) involved choosing priors $p(x)$, or,

equivalently, choosing $-\log(p(x)) =: \rho(x)$

$p(x) = c \cdot \exp(-\rho(x))$ "Gibbs distribution form"

↳ Regularization based ("years" is opinion of authors)

Table 3.1: Evolution of priors for images.

Years	Core concept	Formulae for $\rho(\cdot)$
~ 1970	Energy regularization	$\ x\ _2^2$
1975-1985	Spatial smoothness	$\ \mathbf{L}x\ _2^2$ or $\ \mathbf{D}_v x\ _2^2 + \ \mathbf{D}_h x\ _2^2$
1980-1985	Optimally Learned Transform	$\ \mathbf{T}x\ _2^2 = \mathbf{x}^T \mathbf{R}^{-1} \mathbf{x}$ (via PCA)
1980-1990	Weighted smoothness	$\ \mathbf{L}x\ _{\mathbf{W}}^2$
1990-2000	Robust statistics	$\mathbf{1}^T \mu \{\mathbf{L}x\}$ e.g., Huber-Markov
1992-2005	Total-Variation	$\int_{v \in \Omega} \nabla \mathbf{x}(v) dv = \mathbf{1}^T \sqrt{ \mathbf{D}_v x ^2 + \mathbf{D}_h x ^2}$
1987-2005	Other PDE-based options	$\int_{v \in \Omega} g[\nabla \mathbf{x}(v), \nabla^2 \mathbf{x}(v)] dv$
2005-2009	Field-of-Experts	$\sum_k \lambda_k \mathbf{1}^T \mu_k \{\mathbf{L}_k x\}$
1993-2005	Wavelet sparsity	$\ \mathbf{W}x\ _1$
2000-2010	Self-similarity	$\sum_k \sum_{j \in \Omega(k)} d\{\mathbf{R}_k x, \mathbf{R}_j x\}$
2002-2012	Sparsity methods	$\ \alpha\ _0 \text{ s.t. } \mathbf{x} = \mathbf{D}\alpha$
2010-2017	Low-Rank assumption	$\sum_k \ \mathbf{X}_{\Omega(k)}\ _*$

Image Denoising survey (p. 5)

Wednesday, November 6, 2024 6:18 AM

Other classical denoising

- non-local means (Baudet, Coll, Morel CVPR '05)

optimization method approaches ("MAP") were of great interest academically but almost never made it into Photoshop
(slow, parameters to tune)
...though some variants are used in MRI ...

fancy way to average over nearby (or, similar) pixels

- PDE - inspired methods

- BM3D (Dabov et al. '07)

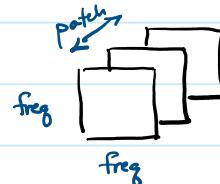
- like JPEG, takes 2D DCT of a patch

Gold standard (prior to deep learning)

- reliable (not too much tuning required)
- fast
- good

- groups similar patches together (cube)

- force joint sparsity on cube



- K-SVD (Elad's): sparsity in a learned "basis"/dictionary

- Improvements were starting to be marginal

'09 paper, "Is Denoising Dead?"

Classical ("mathematical") methods largely ignored color images

"wlog let x be a B+W image..."

RGB, YUV, YCbCr
? ← Chroma
luma

③ Machine Learning Techniques

- Don't explicitly pick a (simple) prior $p(x)$... implicitly learn one from data. i.e. We can learn a denoiser specialized for astronomy, etc.

- Simplest approach: Supervised learning w/ hybrid real/synthetic data

- ML needs a lot of data

- We do have a lot of image data $\{x_k\}_{k=1}^M$ ← "real" data

- To get labeled data, we can artificially add noise

$$y_k = \underbrace{x_k}_{\text{real}} + \underbrace{\nu_k}_{\text{synthetic}} \quad \text{so we have } \{(y_k, x_k)\}_{k=1}^M$$

- Simplest case: $\sigma = \text{constant}$

(see paper for some info on case when we want to do "blind denoising", i.e., σ unknown)

- Define a neural net architecture

- Train using ERM framework (Empirical Risk Minimization) ... as usual

What loss function $\text{dist}(x, \hat{x})$? $\hat{x} = D_\theta(y, \sigma)$

Usually MSE $\|x - \hat{x}\|_2^2$

- Variants:

- σ unknown
→ eg. if you want to tailor to a very specific type of image. "unsupervised"
- No access to clean images $\{x_k\}$
- Beyond AWGN
- Boosting (transfer learning) to specialized datasets (MRI, astronomy, cartoon...)
- SSIM loss function ~ Target perceptual loss

Image Denoising survey (p. 6)

Wednesday, November 6, 2024 6:19 AM

- Some ML architectures are inspired by classical methods
 - Elad calls it "unrolling", slightly similar but not the same as the unrolling we'll discuss when $A \neq I$ "LISTA" Gregor, LeCun ICML '10
 - "Deep K-SVD" and "LIDIA" (like BM3D), see §5
- "Deep image prior": CNN trained on a single image to fit itself (like INR a little)
CNN architecture implicitly regularizes.
Do early stopping! (some disadvantages as classical Lucy-Richardson deblurring!)

④ Beyond Image Denoising The exciting part

- (A) Denoisers ($A = I$) can directly help general inverse problems ($A \neq I$)
 $y = Ax + v$
- (B) Denoisers can directly help generative models to synthesize new images (with applications to (A))
- (C) we can go beyond "MAP" and "MMSE" to get a denoised image w/ much better perceptual quality

(A) $y = Ax + v$, A linear Ex: deblurring, inpainting, demosaicing, Superresolution, tomography (CT, MRI, ...)

$$\text{MAP } \hat{x}_{\text{MAP}} = \arg \min_x \frac{1}{2} \frac{1}{\sigma^2} \|Ax - y\|^2 + c\rho(x) \quad \rho(x) = -\log(p(x)) \text{ as usual}$$

Plug-n-Play idea (Venkatakrishnan, Bouman, Wohlberg, '13 ... though I think it goes back earlier)
Suppose $\rho(x) = \|x\|_1$.

How would people solve the optimization problem?

Two common approaches: ① ADMM or ② proximal gradient descent

both work equally well to motivate key idea.

Paper uses ①, I'll use ②

Image Denoising survey (p. 7)

Wednesday, November 6, 2024

7:04 AM

"Proximal Gradient"

Goal: iteratively minimize

$$f(x), \quad \nabla f(x) = \frac{1}{\sigma^2} A^T (Ax - y)$$

$$\frac{1}{2\sigma^2} \|Ax - y\|^2 + c \cdot \rho(x) \text{ via:}$$

no approximation

$$x_{k+1} = \arg \min_x f(x_k) + \underbrace{\langle \nabla f(x_k), x - x_k \rangle}_{\text{Taylor Series of } f} + \frac{1}{2\gamma} \|x - x_k\|^2 + c \cdot \rho(x)$$

$$(\text{if } \rho(x) = 0, \dots = x_k - \gamma \cdot \nabla f(x_k))$$

which is gradient descent!

$$\star = \arg \min_x \frac{1}{2} \|x - z_k\|^2 + c \cdot \rho(x) \quad (\text{so again, } \rho = 0, \dots = z_k \text{ clearly})$$

$$= \text{prox}_{c \cdot \rho}(z_k). \quad \text{This generalizes projected gradient descent}$$

Classical choices of ρ

were made such that prox_ρ is easy to compute

Ex: $\rho(x) = \|x\|_1 = \sum_{i=1}^n |x_i|$ is separable so whole

prox problem is separable ($\min_{x_1, x_2} f(x_1) + g(x_2)$)
you can do separately?

Wait a minute... \star is the MAP estimator

for denoising ($A = I$)! We (approximately) linearize, denoise, then repeat.

Insight Rather than choose ρ and then calculate prox_ρ ,

$$= D(\cdot, r)$$

lets choose prox_ρ (our denoiser) and leave ρ implicitly defined!

- Advantage: leverage modern denoisers tailored to realistic / specialized datasets

- Disadvantage: math. In interesting cases, we have no reason to believe it must converge...

Image Denoising survey (p. 8)

Wednesday, November 6, 2024 7:18 AM

RED idea "Regularization by Denoising" (Romano, Elad, Milanfar '17)
Somewhat similar spirit to Plug-n-Play award-winning

RED instead

defines ρ via D in a

$$\text{different way: } \rho(x) := x^T(x - D(x; \sigma))$$

with some motivation.

Under conditions on D (smooth, "passive", local homogeneity)

then this ρ is a convex function!

So we get back nice math, convergence.

implicitly
defined by D
now

$$D(y; \sigma) = \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2\sigma^2} \|x - y\|^2 + \rho(x)$$

Exploring this idea further eg. which denoiser to use? which σ ?

$$x_{MAP} = \operatorname{argmin}_x \frac{1}{2} \|Ax - y\|^2 + \underbrace{\sigma^2 \cdot c \cdot -\log(\rho(x))}_{\text{aka } \rho}$$

if ρ is differentiable,

could solve via gradient descent

$$x_{k+1} = x_k - \gamma (A^T(Ax_k - y) - c \cdot \underbrace{\nabla_x \log(\rho(x))}_{\text{"Score function" in statistics}})$$

Going back to 1960's, it's known optimal MMSE denoiser $E(x|y)$

$$\nabla_y \log(p(y)) = \frac{D(y; \sigma_0) - y}{\sigma_0^2}$$

This fits w/ RED definition (if we had optimal denoiser),
and also suggests using a denoiser trained w/ MSE loss.

[RED Algo.] Parameters: stepsize γ , weight c , noise level σ

$$x_{k+1} = x_k - \gamma (A^T(Ax_k - y) + c \cdot (x_k - D(x_k; \sigma)))$$

Image Denoising survey (p. 9)

Wednesday, November 6, 2024 7:39 AM

(B) Generative AI : image synthesis via denoisers

Idea: learn $G_\theta(z)$ to generate an image $x \sim p(x)$
↑
random input, $z \sim N(0, I)$

like (multivariate) inverse CDF sampling... but we won't have explicit formulas

Ex: Variational Autoencoders (VAE)

Ex: Generative Adversarial Networks (GAN) ← were dominant method
for generative AI in image denoising

Ex: Diffusion Models (now dominant)

"Denoising probabilistic models" Ho, Jain, Abbeel. NeurIPS '20

Sohl-Dickstein et al. ICML '15 or "Score-based generative models"
Song + Ermon, NeurIPS '19

How to sample from $p(x)$? (assuming you know $p(x)$)

ex:

- rejection sampling. Very inefficient, essentially inapplicable in higher dim.
- better: Markov Chain Monte Carlo MCMC

Apply a Markov Chain ($x_{k+1} | x_k \perp x_{k-1}, x_{k-2}, \dots$)

so that its stationary distribution is $p(x)$

Welling, Teh '10
popularized in ML
(for other purposes)

- Our models are based on a variant of MCMC:

Langevin Dynamics originally from physics $z_t \sim N(0, I)$

$$x_{t+1} = x_t + \alpha \nabla_x \log(p(x_t)) + \sqrt{2\alpha} z_t$$

α = stepsize

[one motivation: stochastic differential eq'n
 $dX(t) = -\nabla \log(p(X(t)) dt + \sqrt{2/\beta} dB(t)$
then discrete, "Euler-Maruyama" converges to $X \sim p$]

$B = \text{temperature}$
 $B(t) = \text{Brownian motion}$

(under suitable conditions) this gives samples $x_t \sim p$

$\nabla_x \log(p(x))$ is
(or $\nabla_x p(x)$)
implicitly done
via our denoisers

$$\nabla_x P = 2(x - D(x, \sigma))$$

but need tricks to make it practical!

- we initialize w/ some x_0 , but $p(x_0) = 0$ (or very close)

and $\nabla \log(p(x))$ will explode to $-\infty$ ↗ most matrices are not images!

Solution "blur", i.e. add ridge factor,
equivalent to Tikhonov style ...

- slow convergence

Solution: "annealed Langevin Dynamics" ALD (Song, Ermon '19)

Start w/ large σ , run for a while, then decrease σ, \dots ,
until reach target σ

Image Denoising survey (p. 10)

Wednesday, November 6, 2024

10:25 AM

Algo: ALD

Parameters $\{\sigma_i\}_{i=0}^L, \varepsilon, T$

$x_0 \sim N(0, I)$ // size of (vectorized) image

for $i = 1, \dots, L$

$$\alpha_i = \varepsilon \cdot (\sigma_i / \sigma_L)^2$$

for $t = 1, \dots, T$

$$z_t \sim N(0, I)$$

denoiser

$$x_t \leftarrow x_{t-1} + \alpha_i / \sigma_i^2 (D(x_{t-1}, \sigma_i) - x_{t-1}) + \sqrt{2\alpha_i} z_t$$

$$x_0 \leftarrow x_T$$

end

return x_0

(c) "High perceptual quality image recovery": beyond point estimators

MMSE is $\underset{\hat{x}}{\operatorname{argmin}} E \|x - \hat{x}\|^2$

↑
i.e. over $x \sim p$

"MAP" is Bayesian... but
not fully in the Bayesian
spirit, since it's a point estimator

Issue: loss function is MSE, and, averaging over all possible inputs x

Ex: I want to own a dog (one dog), my partner doesn't.

MMSE solution is to own half a dog

Consequence: MMSE results \hat{x} might not look **realistic**

"perception-distortion tradeoff"
realistic? MMSE

One approach: \hat{x} is a sample from $p(x|y)$

Conceptual shift:
random output

- Can be done via GANs (generator vs. discriminator)

- or via diffusion models

$$\nabla_x \left(\log \left(p(x|y) = \frac{p(y|x)p(x)}{p(y)} \right) \right) = \nabla_x p(y|x_t) + \nabla_x p(x)$$

Bayes' Rule
irrelevant
conditional score function

via Eq. (9.4)
in Elad et al.

usual score fun
(via denoiser)

- now MMSE denoiser used to create realistic denoiser!

See Eq (9.2)
for details...
it's slightly
more complicated.

Image Denoising survey (p. 11)

Wednesday, November 6, 2024 10:48 AM

This high-perceptual quality denoiser can be extended
to solve general inverse problems ($y = Ax + v$) as well.
See §9.2. Basically some SVD tricks w, A,
and properties of multivariate normal distribution.

"SNIPS"
approach

Other approaches exist

Some handle when "A" isn't explicitly defined
("stylization", "JPEG deblocking")

Fancier extensions: text-to-image generation (DALL-E)