

# Denoising images with optimization

D.A. Forsyth

University of Illinois

# Denoising Images using Optimization

This chapter uses a master recipe for denoising. Write  $\mathcal{N}$  for a noisy image, and think of denoising as finding a denoised image  $\mathcal{D}$  that is (a) close to  $\mathcal{N}$  and (b) more like a real image. Write

$$\begin{aligned} C(\mathcal{D}) &= [\text{distance from } \mathcal{D} \text{ to } \mathcal{N}] + [\text{unrealism cost for } \mathcal{D}] \\ &= [\text{data term}] + [\text{penalty term}] \end{aligned}$$

and choose a  $\mathcal{D}$  that minimizes this cost function. Methods differ mainly by the penalty term, which has a significant effect on how hard the optimization problem is. This framework leads to very strong denoising methods, at the cost of solving what can be a nasty optimization problem.

# Denoising by weighted least squares - I

For this Chapter, the data term in the master recipe is

$$\sum_{ij} (\mathcal{D}_{ij} - \mathcal{N}_{ij})^2$$

(the ssd of Section 3.4.2). A good reconstruction could smooth the image over quite long scales in regions where  $\mathcal{C}$  is constant. The reconstruction must preserve edges, so the smoothing would need to be over very short scales at edge points. Ideally, smoothing would be along an edge rather than across it. But  $\mathcal{C}$  isn't known (otherwise there would be nothing to do). All this suggests that the penalty function needs to look at gradients in  $\mathcal{D}$ .

# Denoising by WLS - II

- Straighten
  - noisy image  $\mathbf{n}$  into vector  $\mathbf{n}$
  - reconstructed image  $\mathbf{U}$  into vector  $\mathbf{u}$
  
- Cost becomes

$$[\mathbf{u} - \mathbf{n}]^T [\mathbf{u} - \mathbf{n}] .$$

# Denoising by WLS - III

- Idea:
  - reconstructed image should have large gradients only where there is strong evidence in support
  - (version of "pixels are like their neighbors")
- Strong evidence:
  - Use DOG filters to get smoothed gradient of noisy image
  - Where this is big, gradients in reconstruction should be cheap

# Denoising by WLS - IV

- Differentiation is linear, so can write matrices so that gradient of  $u$  is given by

$$\begin{pmatrix} \mathcal{D}_x \\ \mathcal{D}_y \end{pmatrix} \mathbf{u}$$

- What comes out is stacked x and y derivatives

# Denoising by WLS -V

Now write  $\mathcal{A}_x(\mathbf{n})$ ,  $\mathcal{A}_y(\mathbf{n})$  for diagonal matrices of weights obtained from the original image. Because these matrices are diagonal, think of them as producing pixel by pixel weights on the cost of a derivative in  $\mathcal{D}$ . So at a location where the value of  $\mathcal{A}_y$  is small,  $\mathcal{D}$  could have a large  $y$ -derivative, but at locations where the value is large,  $\mathcal{D}$  must have a small  $y$ -derivative.

- Weights could be

$$\frac{1}{|w_i|^\alpha + \epsilon}$$

- Where  $w_i$  is either  $x$  or  $y$  derivative at  $i$ 'th location,  $\epsilon$  is small

# Denoising by WLS - VI

$$\operatorname{argmin}_{\mathbf{u}} [\mathbf{u} - \mathbf{n}]^T [\mathbf{u} - \mathbf{n}] + \lambda \mathbf{u}^T [\mathcal{D}_x^T \mathcal{A}_x^T \mathcal{A}_x \mathcal{D}_x + \mathcal{D}_y^T \mathcal{A}_y^T \mathcal{A}_y \mathcal{D}_y] \mathbf{u}$$



Be close to noisy image



Have big derivatives only if good evidence

where the first term pushes  $\mathbf{d}$  to be like  $\mathbf{n}$ , the second term controls the derivatives of  $\mathbf{d}$  and  $\lambda$  is some weight balancing the two terms. Write  $\mathcal{L} = [\mathcal{D}_x^T \mathcal{A}_x^T \mathcal{A}_x \mathcal{D}_x + \mathcal{D}_y^T \mathcal{A}_y^T \mathcal{A}_y \mathcal{D}_y]$ ; then solving this problem is a matter of solving

$$\mathcal{F}(\lambda)\mathbf{d} = (\mathcal{I} + \lambda\mathcal{L})\mathbf{d} = \mathbf{n}$$



Original



Noisy version



Reconstructions

0.03

0.01

0.005


0.001

Residuals

# Norms - I

$$\begin{aligned} C(\mathcal{D}) &= [\text{distance from } \mathcal{D} \text{ to } \mathcal{N}] + [\text{unrealism cost for } \mathcal{D}] \\ &= [\text{data term}] + [\text{penalty term}] \end{aligned}$$

Q: how to measure the "size" of the penalty?


$$\underset{\mathbf{u}}{\operatorname{argmin}} \quad [\mathbf{u} - \mathbf{n}]^T [\mathbf{u} - \mathbf{n}] + \lambda \mathbf{u}^T [\mathcal{D}_x^T \mathcal{A}_x^T \mathcal{A}_x \mathcal{D}_x + \mathcal{D}_y^T \mathcal{A}_y^T \mathcal{A}_y \mathcal{D}_y] \mathbf{u}$$

# Norms -II

The  $L_2$  norm, defined by

$$\|\mathbf{v}\|_2 = \sqrt{\mathbf{v}^T \mathbf{v}}.$$

$$\operatorname{argmin}_{\mathbf{u}} [\mathbf{u} - \mathbf{n}]^T [\mathbf{u} - \mathbf{n}] + \lambda \mathbf{u}^T [\mathcal{D}_x^T \mathcal{A}_x^T \mathcal{A}_x \mathcal{D}_x + \mathcal{D}_y^T \mathcal{A}_y^T \mathcal{A}_y \mathcal{D}_y] \mathbf{u}$$

This is squared L2 norm (of what?)



# Norms - III

$$\operatorname{argmin}_{\mathbf{u}} [\mathbf{u} - \mathbf{n}]^T [\mathbf{u} - \mathbf{n}] + \lambda \mathbf{u}^T [\mathcal{D}_x^T \mathcal{A}_x^T \mathcal{A}_x \mathcal{D}_x + \mathcal{D}_y^T \mathcal{A}_y^T \mathcal{A}_y \mathcal{D}_y] \mathbf{u}$$

Weighted least squares penalized the squared L2 norm of the weighted gradient. Generally, a vector with small L2 norm can have many small, but non-zero, elements. This is because the square of a small number is very small, and the sum of many very small numbers is still small. The weights in weighted least squares tend to mitigate this, because small gradients have large penalty weights. **Warning:** It is quite common to refer to the *square* of the L2 norm as the L2 norm. I will try not to do this, because it's wrong, but you'll bump into this in the literature rather often.

# The L1 Norm

An alternative is to penalize the *L1 norm* of the gradients. The L1 norm of a vector  $\mathbf{v}$  is defined by

$$\|\mathbf{v}\|_1 = \sum_i |v_i|.$$

# Behavior of L2 norm

A vector with small L1 norm will tend to have zero elements. You can see this by comparing two cases. Write

$$C_2(\mathbf{u}) = \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \frac{\lambda}{2} \mathbf{u}^T \mathbf{u}$$

and notice that the  $\mathbf{u}$  that minimizes  $C_2(\mathbf{u})$  is

$$\frac{1}{1 + \lambda} \mathbf{g}.$$

Notice – even if lambda is very big, g ISN'T zero  
(it's just small)

# Behavior of L1 norm

Now write

$$C_1(\mathbf{u}) = \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \lambda \|\mathbf{u}\|_1$$

and think about the  $\mathbf{u}$  that minimizes  $C_1(\mathbf{u})$ . The penalty term isn't differentiable, which creates some inconvenience, but it is a sum over elements of  $\mathbf{u}$ . Now consider the  $i$ 'th element of  $\mathbf{u}$ . If  $g_i$  is sufficiently large, then it is easy to show that

$$u_i = \frac{g_i}{1 + \lambda}.$$

Now consider what happens when  $g_i = \lambda$ . If  $u_i = 0$ , then the cost will be  $\lambda^2/2$ , but if  $u_i = \epsilon > 0$  where  $\epsilon$  is small, the cost will be  $(1/2)(\lambda^2 + \epsilon^2)$ . This analysis implies correctly that if  $-\lambda < g_i < \lambda$ ,  $u_i = 0$ . In turn, using an L1 norm as a penalty on the gradients tends to cause the reconstruction to have many zero gradients

The L1 norm encourages  $\mathbf{g}$  to have zeros in it!

# Total variation denoising - I

In *total variation denoising*, the penalty is an L1 norm to the gradient. There are a variety of ways of doing this. In one approach, one seeks

$$\operatorname{argmin}_u \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \lambda [\|\mathcal{D}_x \mathbf{u}\|_1 + \|\mathcal{D}_y \mathbf{u}\|_1].$$

Note this cost function isn't differentiable, but it is convex. The optimization problem for this cost function is well understood, and is relatively easily managed (though beyond our scope). However, you should notice that the penalty encourages zeros in the  $x$  and  $y$  components of the gradient, which isn't necessarily the same as zero gradients. One could get a solution where the zeros in the  $x$  components are not aligned with the zeros of the  $y$  components, so the penalty is biased against some gradient directions but not others.



# Total variation denoising - II

An alternative formulation requires a bit more notation. Write  $d_{x,i}(\mathbf{u})$  for the  $i$ 'th component of  $\mathcal{D}_x \mathbf{u}$ , and so on. Then solve

$$\operatorname{argmin}_u \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \lambda \left[ \sum_i \sqrt{d_{x,i}^2 + d_{y,i}^2} \right]$$

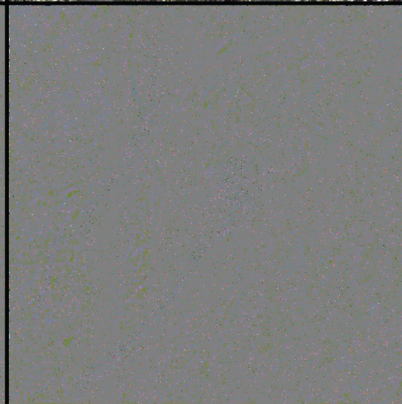
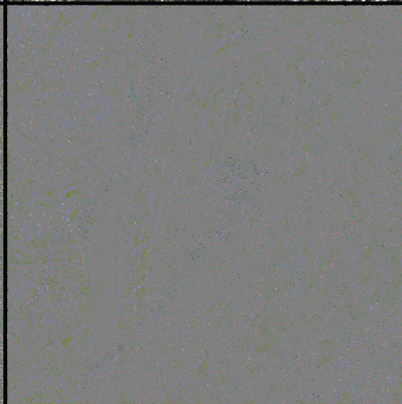
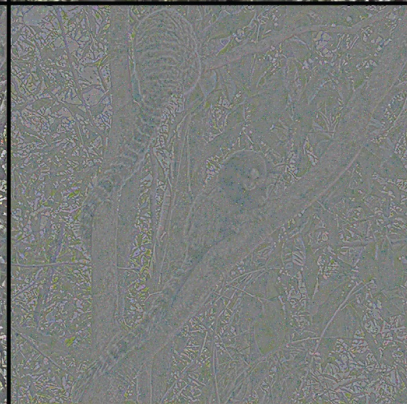
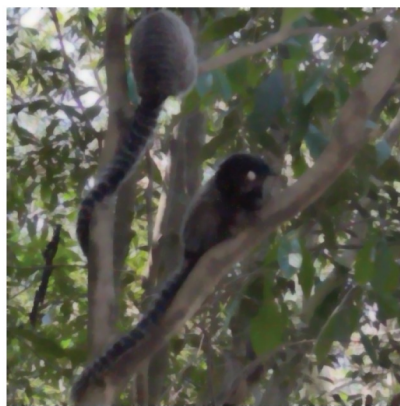
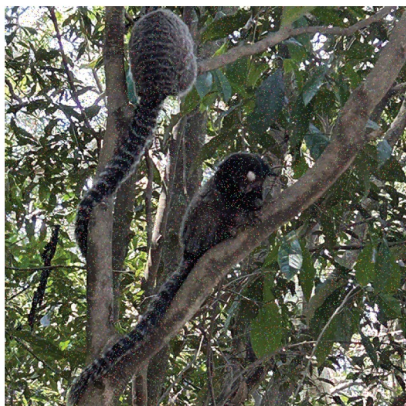
which is also not differentiable. Solutions require rather more elaborate work than solutions for the previous formulation, and tend to be somewhat slower, but are not biased.

100

10

1

0.1



# Deblurring - I

Denoising takes something that isn't quite an image and finds an image that is very like it. Many phenomena can produce something that isn't quite an image. For example, take an image and blur it. The result isn't an image, but it is quite close to one. Recall from Section 41.2 that blurring is a linear operation. Write  $\mathbf{t}$  for the true image in vector form,  $\mathbf{d}$  for the deblurred estimate in vector form,  $\mathbf{b}$  for the observed image in vector form, and  $\mathcal{B}$  for the linear operator that blurs. Assume  $\mathcal{B}$  is known, at least for the moment (**exercises**). Notice  $\mathbf{b}$  is not *exactly* the blurred image. At the very least, there is some error from the numerical representation, and there might be some small noise present, too. Then

$$\mathbf{b} = \mathcal{B}\mathbf{t} + \xi$$

(where  $\xi$  is a vector of very small errors) and least-squares suggests choosing  $\mathbf{d}$  that minimizes

$$(\mathcal{B}\mathbf{d} - \mathbf{b})^T (\mathcal{B}\mathbf{d} - \mathbf{b})$$

which would involve solving

$$\mathcal{B}^T \mathcal{B}\mathbf{d} = \mathcal{B}^T \mathbf{b}.$$

# Deblurring - II

$$\begin{aligned}(\mathcal{B}^T \mathcal{B})^{-1} \mathcal{B}^T \mathbf{b} &= (\mathcal{B}^T \mathcal{B})^{-1} \mathcal{B}^T [\mathcal{B} \mathbf{t} + \boldsymbol{\xi}] \\ &= \mathbf{t} + (\mathcal{B}^T \mathcal{B})^{-1} \mathcal{B}^T \boldsymbol{\xi}.\end{aligned}$$



But this has some really big eigenvalues!

# Regularization

There is a traditional procedure to handle very small eigenvalues in a matrix, known as *regularization*. One seeks a minimum of

$$C(\mathbf{u}) = (\mathcal{B}\mathbf{u} - \mathbf{b})^T (\mathcal{B}\mathbf{u} - \mathbf{b}) + \lambda \mathbf{u}^T \mathbf{u}$$

by solving

$$(\mathcal{B}^T \mathcal{B} + \lambda \mathcal{I})\mathbf{d} = \mathcal{B}^T \mathbf{b}$$

# MLS and TVD

$$\begin{aligned} C(\mathbf{u}) &= [\text{Term comparing } \mathcal{B}\mathbf{u} \text{ to } \mathbf{b}] + [\text{Term evaluating realism of } \mathbf{u}] \\ &= [\text{data term}] + [\text{penalty term}] \end{aligned}$$



Penalty term we used before for MLS or TVD

Blurred input



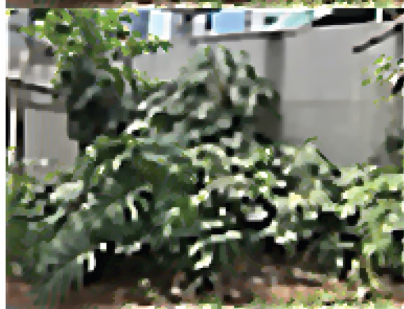
$5e-3$



$5e-2$



$5e-1$



Regularized



$1.5e-5$



$5e-5$



$3e-5$

WLS

Ground truth





Blurred input

$5e-3$



$5e-2$



$5e-1$



Regularized



TV-D

$4e-2$

Ground truth

$6e-2$



$8e-2$



# Think about this....

- 7.5. Section 7.2.6 says that it is harder to deblur a really blurry image than it is to deblur a slightly blurry image, because some eigenvalues of  $(\mathbf{B}^T \mathbf{B})^{-1}$  are very large and very hard to control. Explain.
- 7.6. Assume you know an image is blurred using a gaussian kernel, *and* you know the  $\sigma$  of the kernel. You could deblur using the convolution theorem. What might go wrong if you do?
- 7.7. Assume you know an image is blurred using a gaussian kernel, *and* you know the  $\sigma$  of the kernel. The convolution theorem explains why it is much harder to deblur a heavily blurred image than it is to deblur a lightly blurred image. Explain.
- 7.8. Is the cost function

$$C_1(\mathbf{u}) = \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \lambda \|\mathbf{u}\|_1$$

differentiable?

- 7.9. Imagine you ignore the question of differentiability, and minimize

$$C_1(\mathbf{u}) = \frac{1}{2} [\mathbf{u} - \mathbf{g}]^T [\mathbf{u} - \mathbf{g}] + \lambda \|\mathbf{u}\|_1$$

by gradient descent. You will find the estimated solution you get does not have many zeros in it, even though you are using an L1 norm. Explain what happened.

- 7.10. Section 7.2.4 has: “However, you should notice that the penalty encourages zeros in the  $x$  and  $y$  components of the gradient, which isn’t necessarily the same as zero gradients.” Explain.