

Exercises 8: Spatial smoothing at scale

Jared Fisher

Last Updated November 5, 2016

In this set of exercises, we'll look at smoothing spatial data collected for discrete areal units. In discrete spatial-smoothing problems, "space" is defined not by distance (like in Euclidean space), but by a neighborhood structure that says which areal units are close to which others.

Specifically, let's say that we have observations y_i , each associated with a vertex $s_i \in \mathcal{S}$ in an undirected graph \mathcal{G} with edge set \mathcal{E} . The edge set tells you which sites are neighbors on the graph. For example, in this set of exercises, we'll be looking at data from an fMRI experiment. In this case, s_i is the specific brain region ("voxel"), each y_i is a measurement of the brain activity at that voxel, and \mathcal{E} is the edge structure corresponding to a regular grid.

The underlying statistical model is

$$y_i = x_i + \epsilon_i, \quad i = 1, \dots, n,$$

where x_i is the true denoised signal, and ϵ_i is mean-zero error. Our goal is to estimate x_i in a way that leverages the assumption of spatial smoothness over the underlying graph.

1 Laplacian smoothing

Recall our discussion of Laplacian smoothing from class, which involved the following matrices defined with respect to the graph $\mathcal{G} = (\mathcal{S}, \mathcal{E})$:

- A , the graph adjacency matrix, which is a square $(0, 1)$ -valued matrix with zeros on its diagonal, and a 1 in entry (i, j) if vertices i and j have an edge between them.
- W , the graph degree matrix, which is the diagonal matrix where w_{ii} is the number of neighbors of vertex i .
- L , the Laplacian matrix, defined as $L = W - A$.
- D , the oriented edge matrix of the graph. Letting $m = |\mathcal{E}|$ be the size of the edge set, D is the $m \times n$ matrix defined as follows. If $(j, k), j < k$ is the i th edge in \mathcal{E} , then the i th row of D has a 1 in position j , a -1 in position k , and a 0 everywhere else. Thus the vector Dx encodes the set of pairwise first differences across the edges of the graph.

All of these are sparse matrices—although the graph determines how sparse.

(A) Consider the Laplacian smoothing problem:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|y - x\|_2^2 + \frac{\lambda}{2} x^T L x, \quad (1)$$

where L is the graph Laplacian. Show that the Laplacian matrix has the alternate representation $L = D^T D$, and therefore that penalty term $x^T L x$ can be rewritten as

$$x^T L x = \|Dx\|_2^2.$$

$$\begin{aligned}
[D^T D]_{ij} &= \sum_k D_{ik} * D_{kj} \\
&= \begin{cases} \text{number of edges at vertex } j & i = j \\ -1 & (i, j) \\ 0 & i \neq j; i, j \text{ not connected} \end{cases} \\
&= \begin{cases} \text{number of neighbors of vertex } j & i = j \\ -1 & i, j \text{ share an edge} \\ 0 & i \neq j; i, j \text{ not connected} \end{cases} \\
&= [W - A]_{ij} \\
&= L_{ij}
\end{aligned}$$

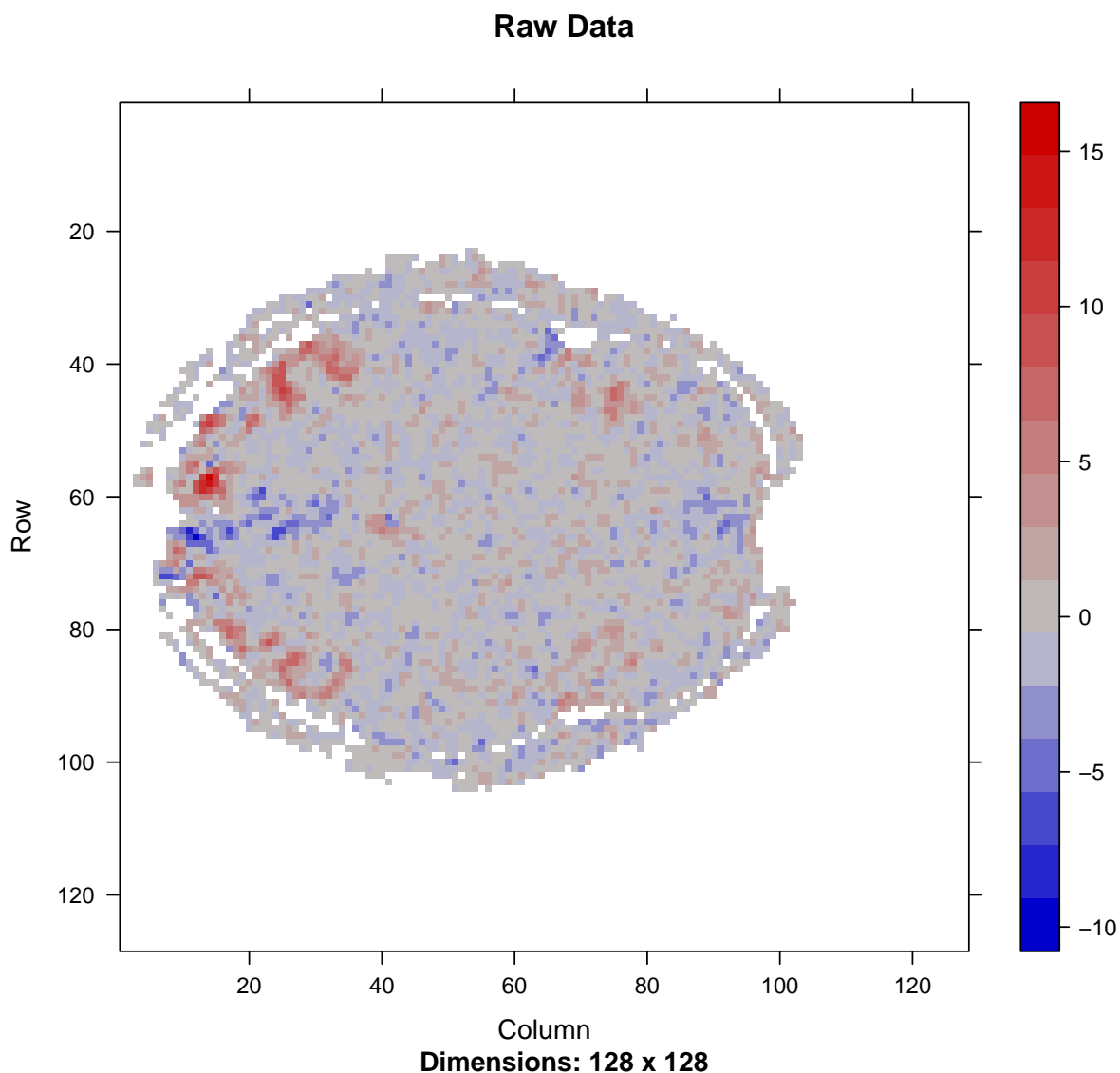
(B) Show that the solution \hat{x} of the Laplacian-smoothing objective solves a linear system

$$C\hat{x} = b,$$

where C and b are matrices you name.

$$\begin{aligned}
\min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} \|y - x\|_2^2 + \frac{\lambda}{2} x^T L x \\
\min_{x \in \mathbb{R}^n} \quad & \frac{1}{2} \|y - x\|_2^2 + \frac{\lambda}{2} \|Dx\|_2^2 \\
\Rightarrow 0 \quad & \stackrel{set}{=} \frac{\partial}{\partial x} \left[\frac{1}{2} \|y - x\|_2^2 + \frac{\lambda}{2} \|Dx\|_2^2 \right] \\
& = -(y - \hat{x}) + \lambda D^T D \hat{x} \\
& = -y + I \hat{x} + \lambda D^T D \hat{x} \\
& = -y + (\lambda D^T D + I) \hat{x} \\
(\lambda D^T D + I) \hat{x} & = y
\end{aligned}$$

(C) Obtain the data “fmri-z.csv” from the class website. This is a 128x128 cross-sectional slice of data from an fMRI experiment on spatial memory, from the laboratory of Russ Poldrack at Stanford University. The data matrix reflects the underlying grid structure of the graph, i.e. entry (i, j) in the matrix corresponds to site (i, j) on the grid. Make a nice heatmap of the data so you can see the spatial structure.



Now consider the following four algorithms for solving the above linear system:

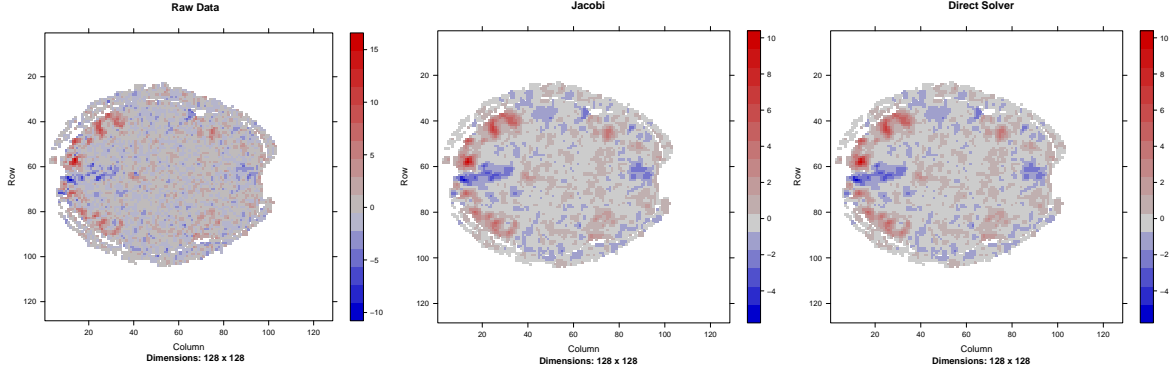
1. a direct solver that uses a sparse matrix factorization (e.g. sparse Cholesky or sparse LU) of the coefficient matrix C .
2. the Gauss-Seidel method.
3. the Jacobi iterative method.
4. the conjugate gradient method. (Specifically, the conjugate gradient method for solving linear systems. There is also a nonlinear conjugate gradient method, which is not what you want.)

You obviously know about the direct solver, but read up on the other three (all of which are iterative methods). Nocedal and Wright have an extensive discussion of the conjugate-gradient method in their *Numerical Optimization* book, while the other two are simple enough to be understood from their Wikipedia pages.

Pick the direct solver (method 1) and any *one* of the other three methods—whichever you think makes the most sense for the problem, in light of having read about them all. Implement both these methods

as *efficiently as you can*, and use them both to solve the Laplacian smoothing problem for our fMRI data. Make sure you get the same results from both methods, up to small numerical errors. Compare them, both in terms of their speed and their ease of implementation. For now, pick λ by eye. Show a nice picture of the raw data side by side with the denoised \hat{x} .

Below we have, in order, the raw data, jacobi-smoothed, and directly-smoothed. Both algorithms finish very quickly, but R's sparse linear system solver takes 2 milliseconds, while mine requires around 480 milliseconds. Any differences between the two methods' results are less than 0.00001. Below we plot them against the original. (Thanks Giorgio for showing us that sparse entries just white themselves out!)



I have given you an R function that constructs the oriented edge matrix D for a two-dimensional grid of size $n_x \times n_y$. You can find it in the file `maked2_sparse.R`.

Also, you can implement all three iterative methods if you want! Just do one iterative method at a minimum.

2 Graph fused lasso

Now consider a version of the spatial smoothing problem where we change the ℓ_2 penalty to an ℓ_1 penalty:

$$\min_{x \in R^n} \frac{1}{2} \|y - x\|_2^2 + \lambda \|Dx\|_1, \quad (2)$$

where we recalled that D is the oriented edge matrix of the graph defined earlier. The penalty term rewards the solution for having small absolute first differences across the edges in the graph. This is called the graph fused lasso, or graph-based total-variation denoising.

Since this doesn't have a closed-form solution, we'll solve this problem via ADMM. There is a fairly obvious ADMM approach, based on rewriting the problem as

$$\begin{aligned} \min_{x \in R^n} \quad & \frac{1}{2} \|y - x\|_2^2 + \lambda \|r\|_1 \\ \text{subject to} \quad & Dx = r. \end{aligned} \quad (3)$$

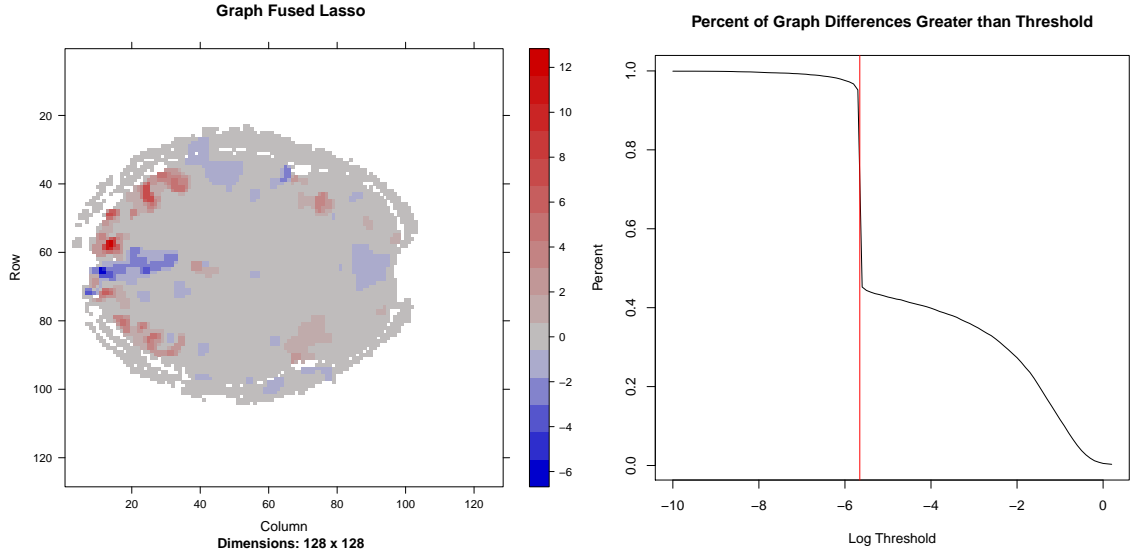
However, this simple approach turns out to be a lot less efficient than the ADMM <https://arxiv.org/abs/1203.1828> described in this paper. (Although I'm biased here, I think Wes Tansey's description in <https://arxiv.org/abs/1411.6144> Section 5 of this paper is much clearer.) The more complex ADMM described in these references leads to more efficient updates that avoid the most expensive matrix operations of the "simple" approach.

Feel free to pick either the "simple and slow" ADMM based on equation (3), or the better one described in either of these links. Implement it, and compare the answer you get on the fMRI data to Laplacian

smoothing. Again, feel free to do something more clever here like cross-validation, but picking λ by eye is fine here.

Note: if you want to get *really* clever, you can optionally use the “proximal-stacking” ADMM algorithm described <https://arxiv.org/abs/1411.0589> in this paper . This will probably require a big investment of time on your part for understanding the approach, but it does lead to a state-of-the-art method for this problem that is *much* faster than either of the algorithms described above. (A related reference is <https://arxiv.org/abs/1505.06475> this technical report; in the case of a grid graph like we have here, the approach described here reduces to the Barbera/Sra technique. I’m again biased here, but I think this paper is easier to understand.)

Here I use the ADMM method used by Tansey and Scott. As presently coded, it is slower than the algorithms used above, taking about 2 minutes of run time. Below is the heatmap of from this optimization routine, next to a “differences” plot. The heatmap shows clearly defined boundaries between brain sections and a smaller number of colored regions than his brother plots, which we would expect to see with the ℓ_1 penalty. The differences plot accompanying it shows the how similar, or not, this smoothed graph is to the others. I chose the direct solver for the comparison. While I’m not definite, I believe we see the effect of the values shrunk to zero around $\exp(-5.5)$. Sure enough, when looking at the histogram/density plot of the errors, it is essentially spike and slab, with a spike at ≈ -0.0035 , and $\log(0.0035) = -5.65$



3 Further notes from James

Remember the algebraic multigrid approach that allows conjugate gradient to fit a linear system in almost linear time.

3.1 How to choose λ

Leave-one-out Lemma: assumption $\hat{y} = Sy$, where S is smoothing matrix (i.e. linear smoother)... See Hastie/Tibsharani/Freedman

3.1.1 Degrees of Freedom

Assume two methods for modeling

1. Fit to p variables.
2. Choose p from $D > p$ variables. Then find best fit for these p variables.

Degrees of freedom of the model- modified definition. Suppose we have $\hat{y} = g(y)$.

$$df(\hat{y}) = \frac{1}{\sigma^2} \sum_{i=1}^n cov(\hat{y}_i, y_i),$$

where $\sigma^2 = var(y_i)$. This means that if $\hat{y} = Sy$, then $df(\hat{y}) = trace(S)$. Then for linear regression, where $\hat{y} = X(X^T X)^{-1} X^T y$.

There's a theorem that says that the trace of a projection matrix is equal to its rank. Hence

$$\begin{aligned} df(\hat{y}) &= trace(X(X^T X)^{-1} X^T) \\ &= rank(X(X^T X)^{-1} X^T) \\ &= p \end{aligned}$$

OR

$$\begin{aligned} df(\hat{y}) &= trace(X(X^T X)^{-1} X^T) \\ &= trace((X^T X)^{-1} X^T X) \\ &= trace(I_p) \\ &= p. \end{aligned}$$

3.1.2 Stein's Lemma

shows an easier way to calculate degrees of freedom than the direct definition, e.g. used to calculate degrees of freedom for lasso.

3.2 Graph Fused Lasso idea

Modified ADMM: $Dx = r$,

$$\begin{aligned} \arg \min_x \quad & \frac{1}{2} \|y - x\|_2^2 + \lambda \|r\|_1 + I_\infty(x, r) \\ \text{s.t.} \quad & x = z, r = s \end{aligned}$$