# Graph-based Penalized Fitting/Smoothing and Semi-supervised Learning

Stephen Vardeman

Analytics Iowa LLC

ISU Statistics and IMSE

# Semi-supervised learning

Consider $N$ complete data cases $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N)$ and $M \geq 0$ additional data cases where only inputs $\mathbf{x}_{N+1}, \ldots, \mathbf{x}_{N+M}$ are available. There is no necessity here that $M > 0$, but it can be so in the event that predictions are desired at $\mathbf{x}_{N+1}, \ldots, \mathbf{x}_{N+M}$ whose values might not be in the training set. Where there are $M > 0$ genuine "unlabeled cases" whose inputs are assumed to come from the same mechanism as the inputs $\mathbf{x}_1, \ldots, \mathbf{x}_N$ and might be used to more or less "fill in" the relevant part of the input space not covered by the complete/labeled data cases, the terminology **semi-supervised learning** is sometimes used to describe the building of a predictor for $y$ at all $N + M$ input vectors. The case $M = 1$ might be used to simply make a single prediction at a single input not exactly seen in a "usual" training set of $N$ complete data pairs.

# Set-up

Suppose that following the development of Section 2.4.3, one can make an adjacency matrix based on the $N + M$ input vectors,

$$\mathbf{S} = (s_{ij})_{\substack{i=1,\ldots,N+M \\ j=1,\ldots,N+M}} = \begin{pmatrix} \mathbf{S}_\mathsf{L} & \mathbf{S}_\mathsf{LU} \\ N \times N & N \times M \\ \mathbf{S}_\mathsf{UL} & \mathbf{S}_\mathsf{U} \\ M \times N & M \times M \end{pmatrix}$$

and corresponding Laplacian and symmetric normalized Laplacian matrices, respectively

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_\mathsf{L} & \mathbf{L}_\mathsf{LU} \\ N \times N & N \times M \\ \mathbf{L}_\mathsf{UL} & \mathbf{L}_\mathsf{U} \\ M \times N & M \times M \end{pmatrix} \text{ and } \mathbf{L}^* = \begin{pmatrix} \mathbf{L}_\mathsf{L}^* & \mathbf{L}_\mathsf{LU}^* \\ N \times N & N \times M \\ \mathbf{L}_\mathsf{UL}^* & \mathbf{L}_\mathsf{U}^* \\ M \times N & M \times M \end{pmatrix}$$

# Graph-penalized fitting/smoothing

With

$$\mathbf{Y}_{(N+M)\times 1} = \begin{pmatrix} \mathbf{Y}_{L} \\ {\scriptstyle N\times 1} \\ \mathbf{Y}_{U} \\ {\scriptstyle M\times 1} \end{pmatrix}$$

one might wish to produce a vector of smoothed/fitted values $\hat{\mathbf{Y}}_{(N+M)\times 1}$ such that entries corresponding to input vectors with large adjacencies tend to be alike. This is possible in way highly reminiscent of the material in Sections 5.1 and 5.3. For $\mathbf{v} \in \Re^{N+M}$ written as

$$\mathbf{v}_{(N+M)\times 1} = \begin{pmatrix} \mathbf{v}_{L} \\ {\scriptstyle N\times 1} \\ \mathbf{v}_{U} \\ {\scriptstyle M\times 1} \end{pmatrix}$$

consider the optimization problem in $\Re^{N+M}$

$$\underset{\mathbf{v}\in\Re^{N+M}}{\text{minimize}} \left( (\mathbf{Y}_{L} - \mathbf{v}_{L})' (\mathbf{Y}_{L} - \mathbf{v}_{L}) + \lambda \mathbf{v}'\mathbf{L}\mathbf{v} \right) \tag{1}$$

# Effects of graph-penalized fitting

The developments of Section 2.4.3 show that upon expanding $\mathbf{v}$ in terms of the $N + M$ (orthonormal) eigenvectors of $\mathbf{L}$ (or $\mathbf{L}^*$) it follows that components of $\mathbf{v}$ that are multiples of late eigenvectors (ones with small eigenvalues)

1. have similar entries for cases with large adjacencies, and

2. are relatively lightly penalized in the minimization.

This strongly suggests that solutions to the optimization problem (1) will provide smoothed prediction vectors $\hat{\mathbf{Y}}$ where entries with corresponding inputs with large adjacencies are similar.

# Ryan and Culp results

Recent work of Culp and Ryan provides theory, methods, and software for solving the problem (1) and many nice generalizations of it (including consideration of losses other than SEL that produce methods for classification problems). Here we will provide (only) the explicit solution that is available for the SEL problem.

It turns out that the problem (1) and generalizations of it separate nicely into two parts. That is

$$\hat{\mathbf{Y}}_U^{opt} = -\mathbf{L}_U^{-1}\mathbf{L}_{UL}\hat{\mathbf{Y}}_L^{opt} \tag{2}$$

(or the same with $\mathbf{L}^*$s replacing $\mathbf{L}$s) where $\hat{\mathbf{Y}}_L^{opt} = \mathbf{v}_L$ solving

$$\underset{\mathbf{v}_L \in \Re^N}{\text{minimize}} \left( (\mathbf{Y}_L - \mathbf{v}_L)' (\mathbf{Y}_L - \mathbf{v}_L) + \lambda \mathbf{v}_L' \tilde{\mathbf{L}}_L \mathbf{v}_L \right) \tag{3}$$

for $\tilde{\mathbf{L}}_L = \mathbf{L}_L - \mathbf{L}_{LU}\mathbf{L}_U^{-1}\mathbf{L}_{UL}$ (or, again, the same with $\mathbf{L}^*$s replacing $\mathbf{L}$s).

# SEL solution

The matrix optimization problem (3) is familiar and its solution a simple consequence of vector calculus

$$\hat{\mathbf{Y}}_L^{\text{opt}} = \left(\mathbf{I} + \lambda \tilde{\mathbf{L}}_L\right)^{-1} \mathbf{Y}_L$$

This is exactly parallel to the results in Sections 5.1 and 5.3. $\left(\mathbf{I} + \lambda \tilde{\mathbf{L}}_L\right)^{-1}$ (and its starred version) is a smoother/shrinker matrix. (Further, the matrix $-\mathbf{L}_U^{-1}\mathbf{L}_{UL}$ in display (2) and its starred version are stochastic matrices and entries of $\hat{\mathbf{Y}}_U^{\text{opt}}$ are averages of the elements of $\hat{\mathbf{Y}}_L^{\text{opt}}$.) Presumably, in an application $\lambda$ can be chosen by cross-validation using the labeled cases to make a CVMSPE.