

Least absolute deviations

Zé Vinícius

July 14, 2018

1 Least absolute deviations

ℓ_1 -norm optimization has been effectively used to estimate models on non-Gaussian noise, especially noise whose distributions possess heavy tails, such as Laplacian noise.

In this note, I followed the Majorization-Minimization (MM) framework in to derive the maximum a posteriori (MAP) affine model subject to Laplacian noise and with prior information that the model coefficients follow a joint Laplacian distribution.

This problem is classically known as the Least Absolute Deviations and can be solved by a class of algorithms known as Iteratively Reweighted Least-Squares.

Mathematically, the cost function is given as follows

$$L(\boldsymbol{\beta}) = \|\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\|_1 + \lambda\|\boldsymbol{\beta}\|_1, \quad (1)$$

in which \mathbf{y} is the $n \times 1$ data vector, \mathbf{X} is the $n \times m$ design matrix, $\boldsymbol{\beta}$ is the $m \times 1$ vector of parameters, $\boldsymbol{\Sigma}$ is a diagonal matrix containing the uncertainties of \mathbf{y} , and λ is a real-valued hyperparameter that controls the strength of the regularization. Our goal is to solve $\operatorname{argmin}_{\boldsymbol{\beta}} L(\boldsymbol{\beta})$.

Following the MM framework, we need to find an upper bound for L , let's call it g , i.e., $g(\mathbf{x}|\mathbf{z}) \geq L(\mathbf{x})$, equality holding at $\mathbf{x} = \mathbf{z}$. Additionally, we hope that g will be "nicier" than L , meaning that we expect to be easier to find a point \mathbf{z}_{t+1} , such that $g(\mathbf{z}_{t+1}|\mathbf{z}_t) \leq g(\mathbf{z}_t|\mathbf{z}_t)$. The best scenario being that \mathbf{z}_{t+1} is a minimizer of $g(\cdot|\mathbf{z}_t)$. If we can construct such g , we can generate a sequence of points that will lead to a feasible point of L under some mild conditions¹.

¹The interested reader is suggested to check out reference (Sun et. al. 2016) for details.

Using eq. (16) from (Sun et. al. 2016), it follows that one possible g can be constructed as

$$g(\beta_{k+1}|\beta_k) = \frac{1}{2} \left\{ \|\Sigma'_k{}^{-1}(\mathbf{y} - \mathbf{X}\beta_{k+1})\|_2^2 + \|\Sigma'_k{}^{-1}(\mathbf{y} - \mathbf{X}\beta_k)\|_1 + \lambda \left[\frac{\|\beta_{k+1}\|_2^2}{\|\beta_k\|_1} + \|\beta_k\|_1 \right] \right\}, \quad (2)$$

in which $\Sigma'_k = \text{diag}(\mathbf{w})$, $w_i = \sigma_i|y_i - \mathbf{X}_i\beta_k|$.

It can be noticed that

$$\text{argmin}_{\beta_{k+1}} g(\beta_{k+1}|\beta_k) = \text{argmin}_{\beta_{k+1}} \|\Sigma'_k{}^{-1}(\mathbf{y} - \mathbf{X}\beta_{k+1})\|_2^2 + \lambda \frac{\|\beta_{k+1}\|_2^2}{\|\beta_k\|_1} \quad (3)$$

which is the least squares cost function with a ℓ_2 regularization component, which has an analytical minimum given as

$$\beta_{k+1} = \left(\mathbf{X}^T \Sigma'_k{}^{-1} \mathbf{X} + \frac{\lambda}{\|\beta_k\|_1} \mathbf{I} \right)^{-1} \mathbf{X}^T \Sigma'_k{}^{-1} \mathbf{y}. \quad (4)$$

References

1. Sun, Y *et. al.*, Majorization-Minimization Algorithms in Signal Processing, Communications, and Machine Learning. *IEEE Transactions on Signal Processing*, 2016.