

# Least-squares parameter estimation in the presence of outliers

Dmitry I. Kabanov

May 13, 2019

Least-squares algorithm is a widely used tool for estimating parameters of a functional relation generated by collected data. The basic formulation is the following. We are given data  $\mathbf{D} = \{D_k\} = \{x_k, Y_k\}$ ,  $k = 1, \dots, N$  with the corresponding errors  $\{\epsilon_k\}$  for the observations  $Y_k$ . Besides, there is an assumption on the functional relation that maps independent variables  $\{x_k\}$  to the noiseless data  $\{F_k\}$ :  $F_k = f(\{x_k\}, \mathbf{X})$ , where  $\mathbf{X} \in \mathbb{R}^M$  is the vector of the model parameters. The aim is to estimate  $\mathbf{X}$  such that the observations  $\{Y_k\}$  fit the model as close as possible.

In the framework of the Bayesian data analysis, we can write that the distribution of  $\mathbf{X}$  can be related to the distribution of the data:

$$\text{prob}(\mathbf{X}|\mathbf{D}) \propto \text{prob}(\mathbf{D}|\mathbf{X}) \times \text{prob}(\mathbf{D}), \quad (1)$$

which, by taking assumption that all the data are independent of each other and that the observational noise obeys to the Gaussian distribution, reduces the problem of finding best estimate of  $\mathbf{X}$  to finding a minimum of the function

$$L = \text{constant} - \sum_{k=1}^N R_k^2, \quad (2)$$

where  $R_k = (F_k - Y_k) / \sigma_k$  are the residuals.

The above formulation of the least-squares algorithm is the most basic one, in which all observations  $\{Y_k\}$  are treated equally. However, the following question arises: what if some of the observations are outliers, that is, erroneous measurements that do not follow the same pattern as other observations?

In this project we consider an extension of the above procedure to accommodate the presence of outliers. To conduct this, we extend Eq. (1) by adding extra parameter  $\sigma$  which represents that observation error over some threshold  $\sigma_0$ ; then the likelihood function in Eq. (1) becomes the marginal PDF. Consequently, the right-hand side of Eq. (2) will change, and the new formulation will allow to estimate parameters  $\mathbf{X}$  more robustly.

We consider two distributions of the error  $\sigma$ :

$$\text{prob}_1(\sigma|\sigma_0) = \sigma_0 / \sigma^2 \text{ for } \sigma \geq \sigma_0 \text{ and } 0 \text{ otherwise}$$

and

$$\text{prob}_2(\sigma|\sigma_0) = \beta \delta(\sigma - \gamma \sigma_0) + (1 - \beta) \delta(\sigma - \sigma_0),$$

where  $0 \leq \beta \ll 1$  and  $\gamma \gg 1$ , and  $\delta$  is the Dirac delta function.

We apply all three formulations of the least-squares procedure to synthetic data with added small-amplitude Gaussian noise, and several data points being corrupted by large-amplitude noise. To assess the performance of different formulations, we compare obtained parameter estimates and their uncertainties.