

# Maximum Likelihood, Robust Linear Regression

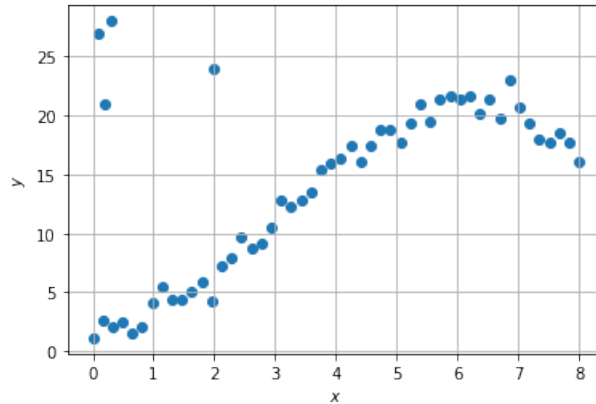
Submit a PDF of your answers to Canvas.

1. A common approach to regression is to imagine the observed data was generated by an additive noise model:

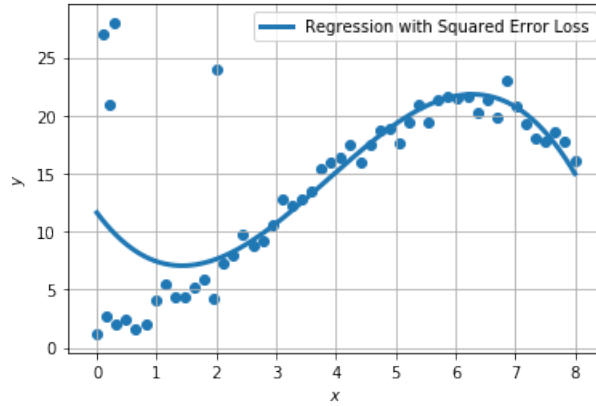
$$y = \mathbf{x}^T \mathbf{w} + Z$$

where  $Z$  is additive noise,  $\mathbf{x}$  is a feature vector and  $\mathbf{w}$  is an unknown vector of coefficients. When  $Z \sim \mathcal{N}(0, \sigma^2)$ , we saw that maximum likelihood estimation for  $\mathbf{w}$  was equivalent to minimizing the *square error loss*. One shortcoming of a squared error loss function is that it results in a fit that is not robust to outliers, since the loss function is (of course) proportional to the *squared* distance between the true value and the prediction.

The scatter plot below shows a dataset  $\mathcal{D} = \{(y_i, \mathbf{x}_i)\}_{i=1}^n$  that includes 4 outliers. The outliers are in the upper left of the plot.



Our goal is to find a *robust* degree-3 polynomial that fits the data points. Recall that we can fit a 3-degree polynomial by first creating a feature vector  $\mathbf{x} = [x^3 \ x^2 \ x^1 \ 1]^T$  for each data point, stacking these as the rows of a  $n \times 4$  matrix  $\mathbf{X}$ , and computing  $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  (which is the closed form for minimizing the squared error loss function). The resulting curve is shown on the plot below. Note the impact of the outliers.



In this problem, to make our fit more robust to outliers, we will assume that the additive noise follows a Laplace distribution. More precisely, we will assume the data was generated by

$$y = \mathbf{w}^T \mathbf{x} + Z$$

where  $Z \sim \text{Lap}(0, 1)$ . Recall the standard Laplace distribution  $\text{Lap}(0, 1)$ :

$$p(z) = \frac{1}{2} e^{-|z|}$$

- Find an expression for the distribution of  $y$  for a fixed  $\mathbf{x}$ ,  $\mathbf{w}$ , denoted  $p_{\mathbf{w}}(y)$ .
- You collect independent measurements and store them in a dataset  $\mathcal{D} = \{(y_i, \mathbf{x}_i)\}_{i=1}^n$ . Find an expression for the distribution of  $y_1, y_2, \dots$  for a fixed  $\mathbf{w}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , denoted  $p_{\mathbf{w}}(\mathbf{y})$ .
- Recall that the maximum likelihood estimate of  $\mathbf{w}$  given  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$  is given by

$$\arg \max_{\mathbf{w}} p_{\mathbf{w}}(\mathbf{y})$$

which is equivalent to

$$\arg \max_{\mathbf{w}} \log(p_{\mathbf{w}}(\mathbf{y})).$$

Use your result for  $p_{\mathbf{w}}(\mathbf{y})$  to simplify the expression above as much as possible. Make use of the notation for the  $\ell_1$  norm:  $\sum |a_i| = \|\mathbf{a}\|_1$ .

- Why is the additive Laplacian noise model more robust to outliers?

**SOLUTION:**

- a) Adding a constant to a pdf shifts the pdf, and  $p_w(y) = \frac{1}{2}e^{-|y-\mathbf{w}^T \mathbf{x}|}$ .
- b) Since the measurements are independent,  $p_w(\mathbf{y}) = \prod_{i=1}^n \frac{1}{2}e^{-|y_i-\mathbf{w}^T \mathbf{x}_i|}$ .
- c) Note that

$$\begin{aligned} \arg \max_{\mathbf{w}} \log p_w(\mathbf{y}) &= \arg \max_{\mathbf{w}} \sum -\log(2) + \sum -|y - \mathbf{w}^T \mathbf{x}_i| \\ &= \arg \min_{\mathbf{w}} \sum |y_i - \mathbf{w}^T \mathbf{x}_i| \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_1. \end{aligned}$$

The results of the additive Gaussian model and the Laplace model are shown below. Note that unlike the least squares estimate,  $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , which has a closed form, we have to use numerical optimization to compute the maximum likelihood estimate of  $\mathbf{w}$  under the Laplacian noise assumption.

