

# Gradient Descent for Logistic-Weighted Kernel-Based Point and Interval Estimation

Will Nickols, Srihari Ganesh

August 7, 2022

## 1 Model

Let  $y$  be the true value for the continuous condition we are interested in. We want to create a probability density function  $f(\hat{y})$  that models the probability of the true crystallization condition being equal to some potential  $\hat{y}$ . We want to maximize the probability assigned to some small interval around the value of the true condition,  $\int_{y(1-\delta)}^{y(1+\delta)} f(\hat{y})d\hat{y}$  for some small  $\delta$ , or, equivalently, we want to minimize the area of the fit density function that falls outside of that interval,  $1 - \int_{y(1-\delta)}^{y(1+\delta)} f(\hat{y})d\hat{y}$ . Because many of the conditions are right skewed, we choose an interval with bounds multiplicatively rather than additively dependent on  $y$ . This probability density can be created by applying a Gaussian kernel to a set of known crystallization conditions  $\mathbf{x}$  from similar proteins. With  $x_i$  as the condition of the crystallization condition for protein  $i$ ,  $p_i$  as the Mash p-value of the alignment between protein  $i$  and the target protein,  $h_i$  as the bandwidth of the kernel element for protein  $i$ ,  $n_p$  as the total number of proteins with Mash p-values less than  $\tau$ ,  $\bar{x}$  as the average of the crystallization conditions of all proteins excluding the protein of interest, and  $\eta$  as the standard deviation of the crystallization conditions of all proteins excluding the protein of interest, this density function can be written as follows.

$$f(\hat{y}) = \frac{1}{n_p + 1} \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\eta} \exp \left[ -\frac{\left( \frac{\hat{y} - \bar{x}}{\eta} \right)^2}{2} \right] + \frac{n_p}{n_p + 1} \cdot \frac{1}{n_p \sqrt{2\pi}} \sum_i \frac{1}{h_i} \exp \left[ -\frac{\left( \frac{\hat{y} - x_i}{h_i} \right)^2}{2} \right] \text{ for } i \text{ such that } p_i < \tau \quad (1)$$

Intuitively, this is a kernel density estimate combining the distribution of crystallization conditions for all proteins and the distribution of crystallization conditions for only proteins similar to the protein of interest. However, two issues arise: not all of these similar proteins are equally similar, so they should not be weighted equally, and the optimal bandwidth  $h_i$  of each term is unknown. Both of these issues can be solved simultaneously by allowing  $h_i$  to be a function of  $s_i$ , the sequence identity (specifically, 1 minus the Mash distance between protein  $i$  and the protein of interest). This function should be continuous and decreasing on  $[0, 1]$  because more similar proteins should have a smaller bandwidth for their kernels, and the function should have a codomain of  $(0, \infty)$  because all the weights should be positive, but some could be much larger than others. Therefore, the relationship between  $h_i$  and  $s_i$  will be given by

$$h_i = \frac{c\sigma(w_1 s_i + w_0)}{\int_0^1 \sigma(w_1 x + w_0) dx} = \frac{c\sigma(w_1 s_i + w_0)}{\frac{\ln(e^{-w_1} + e^{w_0}) - \ln(1 + e^{w_0})}{w_1} + 1} \quad (2)$$

where  $\sigma$  is the sigmoid function and  $c$  is a scaling value to be fit. Letting  $\mathbf{x}$  be the vector of conditions of the similar proteins and  $\mathbf{s}$  be the vector of sequence identities of the similar proteins, define the loss as

$$L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta) = 1 - \int_{y(1-\delta)}^{y(1+\delta)} f(\hat{y}) d\hat{y} + \beta(\eta - c)^2$$

with  $f$  as defined in equation 1 with its  $h_i$  as defined in equation 2. We choose to regularize  $c$  against  $\eta$  because a naive bandwidth should be about the standard deviation of the whole observed condition range, not 0. In practice, letting  $c$  vary often causes the model to break down because values of  $x_i$  close to  $y$  generate an extremely steep gradient for  $c$ , pushing  $c$  very close to 0. By creating extremely sharp peaks of density at each  $x_i$ , this undermines the effort to create a smooth probability density and makes numerical integration essentially impossible. Thus, we will fix  $c$  at  $\eta$ . While this forces the average bandwidth to be the standard deviation of all values for the condition, the function is capable of becoming much larger near 0 than near 1, and protein similarities near 0 often do not pass the p-value threshold for inclusion. Thus, in practice, bandwidths can become as small as necessary even with a fixed  $c$ . Still, for generality, we will treat  $c$  as a variable.

(In the implementation of this model, all values  $y$  and  $x$  of the crystallization condition will be divided by their mean to account for the considerable differences in scale between conditions while maintaining the general right skewedness and positive values of all conditions. When predicting values or ranges on the original scale, we can simply generate an estimate on this altered scale and multiply by the condition's mean.)

## 2 Gradient

The specified model enables the fitting of three parameters:  $w_0$ ,  $w_1$ , and  $c$ . Let  $h_i$  be as described above. Let  $\sigma_i = \sigma(w_1 s_i + w_0)$ . Let  $U$  be the  $\sigma$  normalization term  $\int_0^1 \sigma(w_1 x + w_0) dx = \frac{\ln(e^{-w_1} + e^{w_0}) - \ln(1 + e^{w_0})}{w_1} + 1$ . Let  $d_i = y - x_i$ . Let  $z_i = \left(\frac{d_i}{h_i}\right)$ . Let  $m = e^{w_0} + e^{-w_1}$ . Applying the chain rule, the sum rule for derivatives, and the fact that for  $\sigma(x)$ ,  $\frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x))$ , we obtain the following:

$$\begin{aligned} \frac{\partial f(\hat{y})}{\partial w_0} &= \sum_i \frac{\partial f(\hat{y})}{\partial h_i} \frac{\partial h_i}{\partial w_0} \\ \frac{\partial f(\hat{y})}{\partial w_1} &= \sum_i \frac{\partial f(\hat{y})}{\partial h_i} \frac{\partial h_i}{\partial w_1} \\ \frac{\partial f(\hat{y})}{\partial c} &= \sum_i \frac{\partial f(\hat{y})}{\partial h_i} \frac{\partial h_i}{\partial c} \\ \frac{\partial f(\hat{y})}{\partial h_i} &= \frac{1}{(n_p + 1)\sqrt{2\pi}} \frac{\exp(-z_i^2/2)(z_i^2 - 1)}{h_i^2} \\ \frac{\partial h_i}{\partial w_0} &= \frac{c\sigma_i(1 - \sigma_i)}{U} - \frac{c\sigma_i(\frac{e^{w_0}}{m} - \frac{e^{w_0}}{1+e^{w_0}})}{w_1 U^2} \\ \frac{\partial h_i}{\partial w_1} &= \frac{s_i c\sigma_i(1 - \sigma_i)}{U} - \frac{c\sigma_i[\frac{-w_1 e^{-w_1}}{e^{-w_1} + e^{w_0}} - (\ln(e^{-w_1} + e^{w_0}) - \ln(1 + e^{w_0}))]}{w_1^2 U^2} \\ \frac{\partial h_i}{\partial c} &= \frac{h_i}{c} \end{aligned} \tag{3}$$

However, we are actually interested in the integrals of these quantities, so applying the Leibniz integral rule gives the following:

$$\begin{aligned}
\frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial w_0} &= - \int_{y(1-\delta)}^{y(1+\delta)} \frac{\partial f(\hat{y})}{\partial w_0} d\hat{y} \\
\frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial w_1} &= - \int_{y(1-\delta)}^{y(1+\delta)} \frac{\partial f(\hat{y})}{\partial w_1} d\hat{y} \\
\frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial c} &= - \int_{y(1-\delta)}^{y(1+\delta)} \frac{\partial f(\hat{y})}{\partial c} d\hat{y} + 2\beta(c - \eta)
\end{aligned} \tag{4}$$

### 3 Stochastic Gradient Descent and Model Updating

Because of the memory requirements involved in manipulating all the amino acid identity scores at once, we will use stochastic gradient descent to pick a protein at random, determine its amino acid identity against all the other proteins, compute its density function, and update the weights according to the loss. With a learning rate  $\alpha$ , the update statements will be as follows:

$$\begin{aligned}
w_0 &\leftarrow w_0 - \alpha \frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial w_0} \\
w_1 &\leftarrow w_1 - \alpha \frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial w_1} \\
c &\leftarrow c - \alpha \frac{\partial L(y, \mathbf{x}, \mathbf{s}, \bar{x}, \eta, \delta, \beta)}{\partial c}
\end{aligned} \tag{5}$$

The partial derivative of the density function with respect to each parameter will be computed exactly using equation 3, but the Leibniz integrals in equation 4 will be approximated with a left Riemann sum and a  $\Delta\hat{y}$  of  $y/100$ .

### 4 Expected value, mode estimation, and confidence intervals

By linearity of expectation, the expectation of  $f$  is simply

$$\frac{1}{n_p + 1} \bar{x} + \frac{n_p}{n_p + 1} \frac{1}{n_p} \sum_i x_i$$

Because we do not need extreme precision, the approximate mode of the distribution can be found by evaluating the PDF from  $\min(\bar{x}, \min(\{x_i | x_i \in \mathbf{x}\}))$  to  $\max(\bar{x}, \max(\{x_i | x_i \in \mathbf{x}\}))$  with a step size of the difference divided by 1,000 and recording the value of the condition at the maximum value of the PDF. Because the density of an individual kernel input decreases on either side of its mean, the mode is guaranteed to be between these bounds. Because numeric integration over a large sum of variables is computationally costly and we do not need extreme precision, we can find an estimated 95% confidence interval from the 2.5<sup>th</sup> percentile to the 97.5<sup>th</sup> percentile of the kernel density by iterating over crystallization condition values from either end of the distribution. For the 2.5<sup>th</sup> percentile, we begin iterating upwards from  $\hat{y} = \max(\min(\text{condition value for all proteins}), \min(\Phi^{-1}(0.025)\eta + \bar{x}, \{\min(\Phi^{-1}(0.025)h_i + x_i) | x_i \in \mathbf{x}\}))$  and taking steps of the same size as for the mode until

$$\frac{1}{n_p + 1} \Phi \left( \frac{\hat{y} - \bar{x}}{\eta} \right) + \frac{n_p}{n_p + 1} \cdot \frac{1}{n_p} \sum_i \Phi \left( \frac{\hat{y} - x_i}{h_i} \right) \geq 0.025$$

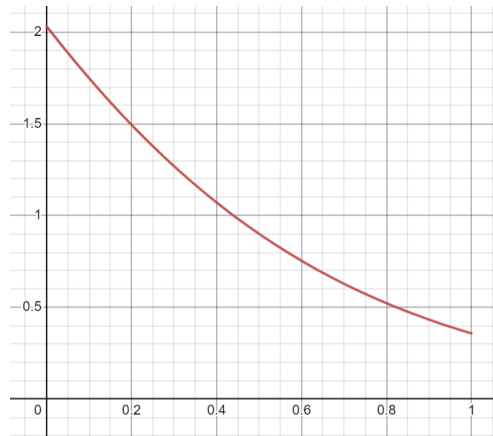
We take this  $\hat{y}$  as the 2.5<sup>th</sup> percentile. Likewise, for the 97.5<sup>th</sup> percentile, we begin iterating downwards from  $\hat{y} = \min(\max(\text{condition value for all proteins}), \max(\Phi^{-1}(0.975)\eta + \bar{x}, \max(\{\Phi^{-1}(0.975)h_i + x_i\} | x_i \in \mathbf{x})))$  and taking steps of the same size until

$$\frac{1}{n_p + 1} \Phi \left( \frac{\hat{y} - \bar{x}}{\eta} \right) + \frac{n_p}{n_p + 1} \cdot \frac{1}{n_p} \sum_i \Phi \left( \frac{\hat{y} - x_i}{h_i} \right) < 0.975$$

We then take the  $\hat{y}$  before the current one (the last one where the expression was greater than 0.975) as the 97.5<sup>th</sup> percentile. As a proof sketch that the minimum of any term's 2.5<sup>th</sup> percentile is less than or equal to the 2.5<sup>th</sup> percentile of the whole kernel density, consider that for any other term besides this minimizer, that term's 2.5<sup>th</sup> percentile must be larger by construction, so it contributes more density to the kernel above its own 2.5<sup>th</sup> percentile and therefore above the minimum 2.5<sup>th</sup> percentile than below either. Therefore, this other term shifts the density of the total kernel density upwards, adding more weight above the minimum 2.5<sup>th</sup> percentile, guaranteeing that the overall 2.5<sup>th</sup> percentile is larger than the minimum of the individual terms' 2.5<sup>th</sup> percentiles. The proof for the 97.5<sup>th</sup> percentile is analogous. Further bounding the search range by the minimum and maximum observed values of the condition ensures that a protein with only a few distantly related proteins doesn't require a massive search space due to very large bandwidths in the kernel density.

## 5 Initialization and runtime

To achieve an initially plausible bandwidth scheme, the following initializations will be chosen:  $w_0 = -1$ ,  $w_1 = -2$ ,  $c = \eta$ . The following image shows the bandwidths produced by these parameters (with  $\eta = 1$ ) with more similar proteins having a lower bandwidth as expected.



The learning rate will start as  $\alpha = 0.1$  and will decay by  $1/(n' \cdot (\text{number of epochs}))$  after each update. If the weights are within a specified convergence radius for a specified number of iterations (0.1 and 20,000 by default), training for that condition's weights will end. Additionally, because the weight updating proceeds

much faster than the Mash distance calculation, the Mash distances will be computed in parallel for a batch of proteins, and the weight updating will proceed sequentially. We fix  $\delta = 0.1$ . Finally, the Mash p-value threshold will be  $\tau = 1/n' \approx 9 \cdot 10^{-6}$ .