

Controlled Average, Variance, and Extrema (CAVE) Functions: Theory, Implementation, and Applications

1 Introduction

We introduce a new group of differentiable functions that can control the average, variance, and extrema (CAVE). With a CAVE function, one can have tight control over certain statistics, namely

$$l = \min(\text{CAVE}(\mathbf{X})) \tag{1}$$

$$h = \max(\text{CAVE}(\mathbf{X})) \tag{2}$$

$$\mu = \mathbb{E}[\text{CAVE}(\mathbf{X})] \tag{3}$$

$$\nu = \text{var}(\text{CAVE}(\mathbf{X})) \tag{4}$$

where l , h , μ , and ν are all chosen by the user. While it is easy to derive functions that control each individually, deriving functions for any combination is not trivial. Here we'll address how each combination can be achieved.

2 Pre-derived Functions

Here we discuss methods that are used to control certain combinations of the minimum, maximum, mean, and variance. Assume that any function $f(\cdot)$ operates element-wise over the input tensor.

2.1 Controlling l and h

The minimum and maximum values are easy to control. ReLU and Softplus for example are unbounded in one direction: any input will have an output in the range $[0, \infty)$. We use functions like these to control the minimum or maximum. More generally, we can denote $f_U(\cdot)$ as any semi-infinitely bound function. Without loss of generality, we'll assume that f_U is lower-bounded by value U . Any function with a finite maximum we can transform by taking the negative.

To limit the output range to our liking, we use f_U as a base function, then appropriately shift it vertically. For setting minimum or maximum values, the CAVE functions are

$$\text{CAVE}(\mathbf{X}; l) = f_U(\mathbf{X}) - U + l \quad (5)$$

$$\text{CAVE}(\mathbf{X}; h) = -(f_U(\mathbf{X}) - U) + h \quad (6)$$

respectively. The new ranges are as desired.

If one wanted to control both the minimum and maximum together, the sigmoid function for example is a common activation function that does just that. Again, we generalize these strictly bounded functions as $f_B(\cdot)$ with lower bound B_l and upper bound B_h . The CAVE function controlling the upper and lower bound is then

$$\text{CAVE}(\mathbf{X}; l, h) = (h - l) \frac{f_B(\mathbf{X}) - B_l}{B_h - B_l} + l. \quad (7)$$

The output range is scaled and shifted to the desired range.

2.2 Controlling μ and ν

We'll use properties of the mean and variance to control these properties. For the mean, we know that

$$\mathbb{E}[\mathbf{X} + b] = \mathbb{E}[\mathbf{X}] + b. \quad (8)$$

It's easy to show that the CAVE function for mean adjustment is then

$$\text{CAVE}(\mathbf{X}; \mu) = \mathbf{X} - \mathbb{E}[\mathbf{X}] + \mu. \quad (9)$$

The variance is likewise easy to derive. We use the property

$$\text{var}(a\mathbf{X}) = a^2 \text{var}(\mathbf{X}) \quad (10)$$

to then derive the CAVE equation

$$\text{CAVE}(\mathbf{X}; \nu) = \sqrt{\frac{\nu}{\text{var}(\mathbf{X})}} \mathbf{X}. \quad (11)$$

Controlling both mean and variance, we can combine the properties to aid in that CAVE derivation:

$$\mathbb{E}[a\mathbf{X} + b] = a\mathbb{E}[\mathbf{X}] + b \quad (12)$$

$$\text{var}(a\mathbf{X} + b) = a^2 \text{var}(\mathbf{X}). \quad (13)$$

Using both properties, the CAVE function to control mean and variance is

$$\text{CAVE}(\mathbf{X}; \mu, \nu) = \sqrt{\frac{\nu}{\text{var}(\mathbf{X})}} (\mathbf{X} - \mathbb{E}[\mathbf{X}]) + \mu. \quad (14)$$

3 Controlling the Remaining Combinations

We previously showed how we can limit the range, but it does not take into account the resulting mean and variance. For example, after applying a sigmoid function, it might be impossible to then apply a linear transform that achieves the desired mean and variance that is guaranteed to be within the desired output range. In order to combat this, we instead first apply a linear transform on the data such that after applying a range limiting function, we achieve all the desired statistics:

$$\text{CAVE}(\mathbf{X}; l, \mu, \nu) = f_U(a\mathbf{X} + b) + l \quad (15)$$

$$\text{CAVE}(\mathbf{X}; h, \mu, \nu) = -f_U(a\mathbf{X} + b) + h \quad (16)$$

$$\text{CAVE}(\mathbf{X}; l, h, \mu, \nu) = (h - l) f_B(a\mathbf{X} + b) + l \quad (17)$$

where a and b are scalars such that the correct statistics are achieved. Note that the mean and variance cannot be simplified in a useful manner in order to solve for a and b . Thus, we solve it numerically.

3.1 CAVE Loss Function

We want to minimize the loss between the calculated mean and/or variance and the desired mean and/or variance by adjusting linear weights a and b . The losses we define are

$$\mathcal{L}(\mathbf{X}, a, b) = \mathcal{L}_\mu(\mathbf{X}, a, b) + \mathcal{L}_\nu(\mathbf{X}, a, b) \quad (18)$$

$$\mathcal{L}_\mu(\mathbf{X}, a, b) = \|E_\mu(\mathbf{X}, a, b)\|_2^2 \quad (19)$$

$$\mathcal{L}_\nu(\mathbf{X}, a, b) = \|E_\nu(\mathbf{X}, a, b)\|_2^2. \quad (20)$$

where

$$E_\mu(\mathbf{X}, a, b) = \mathbb{E}[f(a\mathbf{X}_j + b)] - \mu \quad (21)$$

$$= \frac{1}{N} \sum_{j=1}^N f(a\mathbf{X}_j + b) - \mu \quad (22)$$

$$E_\nu(\mathbf{X}, a, b) = \text{var}(f(a\mathbf{X}_j + b)) - \nu \quad (23)$$

$$= \frac{1}{N} \sum_{j=1}^N f^2(a\mathbf{X}_j + b) - \left(\frac{1}{N} \sum_{j=1}^N f(a\mathbf{X}_j + b) \right)^2 - \nu \quad (24)$$

and $f(\cdot)$ is a range-limited function. Our loss function \mathcal{L} minimizes the l_2 norm of both errors.

Since there are only two variables we are optimizing, we turn to a combination of gradient descent and Newton's method.

3.2 CAVE Optimization

Here we derive how the CAVE function works.

3.2.1 CAVE Loss

We can solve for the gradient analytically by

$$\nabla \mathcal{L}(a, b) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial a} \\ \frac{\partial \mathcal{L}}{\partial b} \end{bmatrix} \quad (25)$$

$$= \begin{bmatrix} \frac{\partial \mathcal{L}_\mu}{\partial a} + \frac{\partial \mathcal{L}_\nu}{\partial a} \\ \frac{\partial \mathcal{L}_\mu}{\partial b} + \frac{\partial \mathcal{L}_\nu}{\partial b} \end{bmatrix} \quad (26)$$

where

$$\frac{\partial \mathcal{L}_\mu}{\partial a} = 2E_\mu \frac{\partial E_\mu}{\partial a} \quad (27)$$

$$\frac{\partial \mathcal{L}_\nu}{\partial a} = 2E_\nu \frac{\partial E_\nu}{\partial a} \quad (28)$$

$$\frac{\partial \mathcal{L}_\mu}{\partial b} = 2E_\mu \frac{\partial E_\mu}{\partial b} \quad (29)$$

$$\frac{\partial \mathcal{L}_\nu}{\partial b} = 2E_\nu \frac{\partial E_\nu}{\partial b}. \quad (30)$$

The Hessian we find by

$$\nabla^2 \mathcal{L}(a, b) = \begin{bmatrix} \frac{\partial^2 \mathcal{L}}{\partial a^2} & \frac{\partial^2 \mathcal{L}}{\partial a \partial b} \\ \frac{\partial^2 \mathcal{L}}{\partial a \partial b} & \frac{\partial^2 \mathcal{L}}{\partial b^2} \end{bmatrix} \quad (31)$$

$$= \begin{bmatrix} \frac{\partial^2 \mathcal{L}_\mu}{\partial a^2} + \frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} & \frac{\partial^2 \mathcal{L}_\mu}{\partial a \partial b} + \frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial b} \\ \frac{\partial^2 \mathcal{L}_\mu}{\partial a \partial b} + \frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial b} & \frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} + \frac{\partial^2 \mathcal{L}_\nu}{\partial b^2} \end{bmatrix} \quad (32)$$

where

$$\frac{\partial^2 \mathcal{L}_\mu}{\partial a^2} = 2 \left(\left(\frac{\partial E_\mu}{\partial a} \right)^2 + E_\mu \frac{\partial^2 E_\mu}{\partial a^2} \right) \quad (33)$$

$$\frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} = 2 \left(\left(\frac{\partial E_\nu}{\partial a} \right)^2 + E_\nu \frac{\partial^2 E_\nu}{\partial a^2} \right) \quad (34)$$

$$\frac{\partial^2 \mathcal{L}_\mu}{\partial a \partial b} = 2 \left(\frac{\partial E_\mu}{\partial a} \frac{\partial E_\mu}{\partial b} + E_\mu \frac{\partial^2 E_\mu}{\partial a \partial b} \right) \quad (35)$$

$$\frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial b} = 2 \left(\frac{\partial E_\nu}{\partial a} \frac{\partial E_\nu}{\partial b} + E_\nu \frac{\partial^2 E_\nu}{\partial a \partial b} \right) \quad (36)$$

$$\frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} = 2 \left(\left(\frac{\partial E_\mu}{\partial b} \right)^2 + E_\mu \frac{\partial^2 E_\mu}{\partial b^2} \right) \quad (37)$$

$$\frac{\partial^2 \mathcal{L}_\nu}{\partial b^2} = 2 \left(\left(\frac{\partial E_\nu}{\partial b} \right)^2 + E_\nu \frac{\partial^2 E_\nu}{\partial b^2} \right). \quad (38)$$

3.2.2 CAVE Error

The mean error function and its derivatives are found as

$$E_\mu = \frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) - \mu \quad (39)$$

$$\frac{\partial E_\mu}{\partial a} = \frac{1}{N} \sum_{j=1}^N \frac{\partial f}{\partial a} \quad (40)$$

$$\frac{\partial E_\mu}{\partial b} = \frac{1}{N} \sum_{j=1}^N \frac{\partial f}{\partial b} \quad (41)$$

$$\frac{\partial^2 E_\mu}{\partial a^2} = \frac{1}{N} \sum_{j=1}^N \frac{\partial^2 f}{\partial a^2} \quad (42)$$

$$\frac{\partial^2 E_\mu}{\partial a \partial b} = \frac{1}{N} \sum_{j=1}^N \frac{\partial^2 f}{\partial a \partial b} \quad (43)$$

$$\frac{\partial^2 E_\mu}{\partial b^2} = \frac{1}{N} \sum_{j=1}^N \frac{\partial^2 f}{\partial b^2}. \quad (44)$$

The variance error function and its derivatives are also found as

$$E_\nu = \frac{1}{N} \sum_{j=1}^N f^2 - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right)^2 - \nu \quad (45)$$

$$\frac{\partial E_\nu}{\partial a} = 2 \left\{ \frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \frac{\partial f}{\partial a} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial E_\mu}{\partial a} \right\} \quad (46)$$

$$\frac{\partial E_\nu}{\partial b} = 2 \left\{ \frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \frac{\partial f}{\partial b} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial E_\mu}{\partial b} \right\} \quad (47)$$

$$\begin{aligned} \frac{\partial^2 E_\nu}{\partial a^2} = 2 \left\{ \frac{1}{N} \sum_{j=1}^N \left[\left(\frac{\partial f}{\partial a} \right)^2 + f(a \mathbf{X}_j + b) \frac{\partial^2 f}{\partial a^2} \right] - \dots \right. \\ \left. \left(\frac{\partial E_\mu}{\partial a} \right)^2 - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^2 E_\mu}{\partial a^2} \right\} \end{aligned} \quad (48)$$

$$\begin{aligned} \frac{\partial^2 E_\nu}{\partial a \partial b} = 2 \left\{ \frac{1}{N} \sum_{j=1}^N \left[\frac{\partial f}{\partial a} \frac{\partial f}{\partial b} + f(a \mathbf{X}_j + b) \frac{\partial^2 f}{\partial a \partial b} \right] - \dots \right. \\ \left. \frac{\partial E_\mu}{\partial a} \frac{\partial E_\mu}{\partial b} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^2 E_\mu}{\partial a \partial b} \right\} \end{aligned} \quad (49)$$

$$\begin{aligned} \frac{\partial^2 E_\nu}{\partial b^2} = 2 \left\{ \frac{1}{N} \sum_{j=1}^N \left[\left(\frac{\partial f}{\partial b} \right)^2 + f(a \mathbf{X}_j + b) \frac{\partial^2 f}{\partial b^2} \right] - \dots \right. \\ \left. \left(\frac{\partial E_\mu}{\partial b} \right)^2 - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^2 E_\mu}{\partial b^2} \right\}. \end{aligned} \quad (50)$$

3.2.3 CAVE Base Function

Lastly, we have the derivatives of the base function f :

$$\frac{\partial f}{\partial a} = \frac{\partial}{\partial a} [f(a \mathbf{X}_j + b)] = f'(a \mathbf{X}_j + b) \mathbf{X}_j \quad (51)$$

$$\frac{\partial f}{\partial b} = \frac{\partial}{\partial b} [f(a \mathbf{X}_j + b)] = f'(a \mathbf{X}_j + b) \quad (52)$$

$$\frac{\partial^2 f}{\partial a^2} = \frac{\partial^2}{\partial a^2} [f(a \mathbf{X}_j + b)] = f''(a \mathbf{X}_j + b) \mathbf{X}_j^2 \quad (53)$$

$$\frac{\partial^2 f}{\partial a \partial b} = \frac{\partial^2}{\partial a \partial b} [f(a \mathbf{X}_j + b)] = f''(a \mathbf{X}_j + b) \mathbf{X}_j \quad (54)$$

$$\frac{\partial^2 f}{\partial b^2} = \frac{\partial^2}{\partial b^2} [f(a \mathbf{X}_j + b)] = f''(a \mathbf{X}_j + b). \quad (55)$$

Given all these values, we can minimize the loss function.

3.3 CAVE Updates

The update function for gradient descent is

$$\begin{bmatrix} a \\ b \end{bmatrix}_{k+1} = \begin{bmatrix} a \\ b \end{bmatrix}_k - \eta \nabla \mathcal{L} \left(\begin{bmatrix} a \\ b \end{bmatrix}_k \right) \quad (56)$$

and the update function for Newton's method is

$$\begin{bmatrix} a \\ b \end{bmatrix}_{k+1} = \begin{bmatrix} a \\ b \end{bmatrix}_k - \eta \left(\nabla^2 \mathcal{L} \left(\begin{bmatrix} a \\ b \end{bmatrix}_k \right) \right)^{-1} \nabla \mathcal{L} \left(\begin{bmatrix} a \\ b \end{bmatrix}_k \right). \quad (57)$$

The above two equations optimize both the mean and variance jointly. If only the mean is specified, the gradient descent step is

$$b_{k+1} = b_k - \eta \frac{\partial \mathcal{L}_\mu}{\partial b} \quad (58)$$

and the Newton step is

$$b_{k+1} = b_k - \eta \left(\frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} \right)^{-1} \frac{\partial \mathcal{L}_\mu}{\partial b}. \quad (59)$$

Similarly, if only the variance is specified, the gradient descent step is

$$a_{k+1} = a_k - \eta \frac{\partial \mathcal{L}_\nu}{\partial a} \quad (60)$$

and the Newton step is

$$a_{k+1} = a_k - \eta \left(\frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} \right)^{-1} \frac{\partial \mathcal{L}_\nu}{\partial a}. \quad (61)$$

To run CAVE functions, we take K_g gradient descent steps to approach the minimum followed by K_n Newton steps to rapidly approach the minimum. Once a and b are found, the data is finally transformed and fed through the base function as $f(t(\mathbf{X}_j, a, b))$.

4 The Gradient of Gradient Descent and Newton's Method

This function, while effective, is memory-intensive when tracking the gradient: many operations require saving the input tensors for backpropagation. Here, we derive the gradient of one step of gradient descent and Newton's method with respect to \mathbf{X}_i given only the input data. The previous section derived how to optimize linear transform variables a and b , but we'll derive how the steps to approach the minimum affect the input data.

Since it is impractical to unravel gradient descent and Newton's method, we'll find the gradient of $\nabla \mathcal{L}(\mathbf{X}, a, b)$ and $(\nabla^2 \mathcal{L}(\mathbf{X}, a, b))^{-1} \nabla \mathcal{L}(\mathbf{X}, a, b)$ both w.r.t. \mathbf{X}_i . These are the values multiplied by the learning rate.

4.1 Gradient of Newton's Method w.r.t. \mathbf{X}_i

We can expand the gradient step by

$$(\nabla^2 \mathcal{L}(\mathbf{X}, a, b))^{-1} \nabla \mathcal{L}(\mathbf{X}, a, b) = \frac{1}{D} \begin{bmatrix} N_a \\ N_b \end{bmatrix}. \quad (62)$$

where

$$N_a(\mathbf{X}, a, b) = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial^2 \mathcal{L}}{\partial b^2} - \frac{\partial \mathcal{L}}{\partial b} \frac{\partial^2 \mathcal{L}}{\partial a \partial b} \quad (63)$$

$$N_b(\mathbf{X}, a, b) = \frac{\partial \mathcal{L}}{\partial b} \frac{\partial^2 \mathcal{L}}{\partial a^2} - \frac{\partial \mathcal{L}}{\partial a} \frac{\partial^2 \mathcal{L}}{\partial a \partial b} \quad (64)$$

$$D(\mathbf{X}, a, b) = \frac{\partial^2 \mathcal{L}}{\partial a^2} \frac{\partial^2 \mathcal{L}}{\partial b^2} - \left(\frac{\partial^2 \mathcal{L}}{\partial a \partial b} \right)^2. \quad (65)$$

We need to find the derivative of each entry w.r.t. each \mathbf{X}_i :

$$\frac{\partial}{\partial \mathbf{X}_i} \left[\frac{N_a(\mathbf{X}, a, b)}{D(\mathbf{X}, a, b)} \right] = \frac{D \frac{\partial N_a}{\partial \mathbf{X}_i} - N_a \frac{\partial D}{\partial \mathbf{X}_i}}{D^2} \quad (66)$$

$$\frac{\partial}{\partial \mathbf{X}_i} \left[\frac{N_b(\mathbf{X}, a, b)}{D(\mathbf{X}, a, b)} \right] = \frac{D \frac{\partial N_b}{\partial \mathbf{X}_i} - N_b \frac{\partial D}{\partial \mathbf{X}_i}}{D^2} \quad (67)$$

where

$$\frac{\partial N_a}{\partial \mathbf{X}_i} = \frac{\partial^2 \mathcal{L}}{\partial a \partial \mathbf{X}_i} \cdot \frac{\partial^2 \mathcal{L}}{\partial b^2} + \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial^3 \mathcal{L}}{\partial b^2 \partial \mathbf{X}_i} - \frac{\partial^2 \mathcal{L}}{\partial b \partial \mathbf{X}_i} \cdot \frac{\partial^2 \mathcal{L}}{\partial a \partial b} - \frac{\partial \mathcal{L}}{\partial b} \cdot \frac{\partial^3 \mathcal{L}}{\partial a \partial b \partial \mathbf{X}_i} \quad (68)$$

$$\frac{\partial N_b}{\partial \mathbf{X}_i} = \frac{\partial^2 \mathcal{L}}{\partial b \partial \mathbf{X}_i} \cdot \frac{\partial^2 \mathcal{L}}{\partial a^2} + \frac{\partial \mathcal{L}}{\partial b} \cdot \frac{\partial^3 \mathcal{L}}{\partial a^2 \partial \mathbf{X}_i} - \frac{\partial^2 \mathcal{L}}{\partial a \partial \mathbf{X}_i} \cdot \frac{\partial^2 \mathcal{L}}{\partial a \partial b} - \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial^3 \mathcal{L}}{\partial a \partial b \partial \mathbf{X}_i} \quad (69)$$

$$\frac{\partial D}{\partial \mathbf{X}_i} = \frac{\partial^3 \mathcal{L}}{\partial a^2 \partial \mathbf{X}_i} \cdot \frac{\partial^2 \mathcal{L}}{\partial b^2} + \frac{\partial^2 \mathcal{L}}{\partial a^2} \cdot \frac{\partial^3 \mathcal{L}}{\partial b^2 \partial \mathbf{X}_i} - 2 \frac{\partial^2 \mathcal{L}}{\partial a \partial b} \cdot \frac{\partial^3 \mathcal{L}}{\partial a \partial b \partial \mathbf{X}_i}. \quad (70)$$

Many of these derivatives have been found in the previous section, but any derivative w.r.t. \mathbf{X}_i has yet to be found. These derivatives of \mathcal{L} are decomposed into their \mathcal{L}_μ and \mathcal{L}_ν counterparts as

$$\frac{\partial^2 \mathcal{L}}{\partial a \partial \mathbf{X}_i} = \frac{\partial^2 \mathcal{L}_\mu}{\partial a \partial \mathbf{X}_i} + \frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial \mathbf{X}_i} \quad (71)$$

$$\frac{\partial^2 \mathcal{L}}{\partial b \partial \mathbf{X}_i} = \frac{\partial^2 \mathcal{L}_\mu}{\partial b \partial \mathbf{X}_i} + \frac{\partial^2 \mathcal{L}_\nu}{\partial b \partial \mathbf{X}_i} \quad (72)$$

$$\frac{\partial^3 \mathcal{L}}{\partial a^2 \partial \mathbf{X}_i} = \frac{\partial^3 \mathcal{L}_\mu}{\partial a^2 \partial \mathbf{X}_i} + \frac{\partial^3 \mathcal{L}_\nu}{\partial a^2 \partial \mathbf{X}_i} \quad (73)$$

$$\frac{\partial^3 \mathcal{L}}{\partial a \partial b \partial \mathbf{X}_i} = \frac{\partial^3 \mathcal{L}_\mu}{\partial a \partial b \partial \mathbf{X}_i} + \frac{\partial^3 \mathcal{L}_\nu}{\partial a \partial b \partial \mathbf{X}_i} \quad (74)$$

$$\frac{\partial^3 \mathcal{L}}{\partial b^2 \partial \mathbf{X}_i} = \frac{\partial^3 \mathcal{L}_\mu}{\partial b^2 \partial \mathbf{X}_i} + \frac{\partial^3 \mathcal{L}_\nu}{\partial b^2 \partial \mathbf{X}_i}. \quad (75)$$

4.1.1 Gradient of Losses w.r.t. \mathbf{X}_i

Using previous definitions, we can find the above loss derivatives. The second order derivatives are

$$\frac{\partial^2 \mathcal{L}_\mu}{\partial a \partial \mathbf{X}_i} = 2 \left(\frac{\partial E_\mu}{\partial \mathbf{X}_i} \cdot \frac{\partial E_\mu}{\partial a} + E_\mu \cdot \frac{\partial^2 E_\mu}{\partial a \partial \mathbf{X}_i} \right) \quad (76)$$

$$\frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial \mathbf{X}_i} = 2 \left(\frac{\partial E_\nu}{\partial \mathbf{X}_i} \cdot \frac{\partial E_\nu}{\partial a} + E_\nu \cdot \frac{\partial^2 E_\nu}{\partial a \partial \mathbf{X}_i} \right) \quad (77)$$

$$\frac{\partial^2 \mathcal{L}_\mu}{\partial b \partial \mathbf{X}_i} = 2 \left(\frac{\partial E_\mu}{\partial \mathbf{X}_i} \cdot \frac{\partial E_\mu}{\partial b} + E_\mu \cdot \frac{\partial^2 E_\mu}{\partial b \partial \mathbf{X}_i} \right) \quad (78)$$

$$\frac{\partial^2 \mathcal{L}_\nu}{\partial b \partial \mathbf{X}_i} = 2 \left(\frac{\partial E_\nu}{\partial \mathbf{X}_i} \cdot \frac{\partial E_\nu}{\partial b} + E_\nu \cdot \frac{\partial^2 E_\nu}{\partial b \partial \mathbf{X}_i} \right) \quad (79)$$

and the third derivatives are

$$\frac{\partial^3 \mathcal{L}_\mu}{\partial a^2 \partial \mathbf{X}_i} = 2 \left(2 \frac{\partial E_\mu}{\partial a} \cdot \frac{\partial^2 E_\mu}{\partial a \partial \mathbf{X}_i} + \frac{\partial E_\mu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\mu}{\partial a^2} + E_\mu \cdot \frac{\partial^3 E_\mu}{\partial a^2 \partial \mathbf{X}_i} \right) \quad (80)$$

$$\frac{\partial^3 \mathcal{L}_\nu}{\partial a^2 \partial \mathbf{X}_i} = 2 \left(2 \frac{\partial E_\nu}{\partial a} \cdot \frac{\partial^2 E_\nu}{\partial a \partial \mathbf{X}_i} + \frac{\partial E_\nu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\nu}{\partial a^2} + E_\nu \cdot \frac{\partial^3 E_\nu}{\partial a^2 \partial \mathbf{X}_i} \right) \quad (81)$$

$$\frac{\partial^3 \mathcal{L}_\mu}{\partial a \partial b \partial \mathbf{X}_i} = 2 \left(\frac{\partial^2 E_\mu}{\partial a \partial \mathbf{X}_i} \cdot \frac{\partial E_\mu}{\partial b} + \frac{\partial E_\mu}{\partial a} \cdot \frac{\partial^2 E_\mu}{\partial b \partial \mathbf{X}_i} + \frac{\partial E_\mu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\mu}{\partial a \partial b} + E_\mu \cdot \frac{\partial^3 E_\mu}{\partial a \partial b \partial \mathbf{X}_i} \right) \quad (82)$$

$$\frac{\partial^3 \mathcal{L}_\nu}{\partial a \partial b \partial \mathbf{X}_i} = 2 \left(\frac{\partial^2 E_\nu}{\partial a \partial \mathbf{X}_i} \cdot \frac{\partial E_\nu}{\partial b} + \frac{\partial E_\nu}{\partial a} \cdot \frac{\partial^2 E_\nu}{\partial b \partial \mathbf{X}_i} + \frac{\partial E_\nu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\nu}{\partial a \partial b} + E_\nu \cdot \frac{\partial^3 E_\nu}{\partial a \partial b \partial \mathbf{X}_i} \right) \quad (83)$$

$$\frac{\partial^3 \mathcal{L}_\mu}{\partial b^2 \partial \mathbf{X}_i} = 2 \left(2 \frac{\partial E_\mu}{\partial b} \cdot \frac{\partial^2 E_\mu}{\partial b \partial \mathbf{X}_i} + \frac{\partial E_\mu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\mu}{\partial b^2} + E_\mu \cdot \frac{\partial^3 E_\mu}{\partial b^2 \partial \mathbf{X}_i} \right) \quad (84)$$

$$\frac{\partial^3 \mathcal{L}_\nu}{\partial b^2 \partial \mathbf{X}_i} = 2 \left(2 \frac{\partial E_\nu}{\partial b} \cdot \frac{\partial^2 E_\nu}{\partial b \partial \mathbf{X}_i} + \frac{\partial E_\nu}{\partial \mathbf{X}_i} \cdot \frac{\partial^2 E_\nu}{\partial b^2} + E_\nu \cdot \frac{\partial^3 E_\nu}{\partial b^2 \partial \mathbf{X}_i} \right). \quad (85)$$

4.1.2 Gradient of Errors w.r.t. \mathbf{X}_i

Similarly, we can find the derivatives of the error functions. The mean error functions are

$$\frac{\partial E_\mu}{\partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial f}{\partial \mathbf{X}_i} \quad (86)$$

$$\frac{\partial^2 E_\mu}{\partial a \partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} \quad (87)$$

$$\frac{\partial^2 E_\mu}{\partial b \partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} \quad (88)$$

$$\frac{\partial^3 E_\mu}{\partial a^2 \partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial^3 f}{\partial a^2 \partial \mathbf{X}_i} \quad (89)$$

$$\frac{\partial^3 E_\mu}{\partial a \partial b \partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial^3 f}{\partial a \partial b \partial \mathbf{X}_i} \quad (90)$$

$$\frac{\partial^3 E_\mu}{\partial b^2 \partial \mathbf{X}_i} = \frac{1}{N} \frac{\partial^3 f}{\partial b^2 \partial \mathbf{X}_i} \quad (91)$$

and the variance error functions are

$$\frac{\partial E_\nu}{\partial \mathbf{X}_i} = \frac{2}{N} \left\{ \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial f}{\partial \mathbf{X}_i} \right\} \quad (92)$$

$$\begin{aligned} \frac{\partial^2 E_\nu}{\partial a \partial \mathbf{X}_i} = \frac{2}{N} \left\{ \frac{\partial f}{\partial a} \frac{\partial f}{\partial \mathbf{X}_i} + f(a \mathbf{X}_i + b) \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} - \dots \right. \\ \left. \frac{\partial E_\mu}{\partial a} \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} \right\} \end{aligned} \quad (93)$$

$$\begin{aligned} \frac{\partial^2 E_\nu}{\partial b \partial \mathbf{X}_i} = \frac{2}{N} \left\{ \frac{\partial f}{\partial b} \frac{\partial f}{\partial \mathbf{X}_i} + f(a \mathbf{X}_i + b) \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} - \dots \right. \\ \left. \frac{\partial E_\mu}{\partial b} \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} \right\} \end{aligned} \quad (94)$$

$$\begin{aligned} \frac{\partial^3 E_\nu}{\partial a^2 \partial \mathbf{X}_i} = \frac{2}{N} \left\{ 2 \frac{\partial f}{\partial a} \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} + \frac{\partial f}{\partial \mathbf{X}_i} \frac{\partial^2 f}{\partial a^2} + f(a \mathbf{X}_j + b) \frac{\partial^3 f}{\partial a^2 \partial \mathbf{X}_i} - \dots \right. \\ \left. 2 \frac{\partial E_\mu}{\partial a} \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} - \frac{\partial^2 E_\mu}{\partial a^2} \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^3 f}{\partial a^2 \partial \mathbf{X}_i} \right\} \end{aligned} \quad (95)$$

$$\begin{aligned} \frac{\partial^3 E_\nu}{\partial a \partial b \partial \mathbf{X}_i} = \frac{2}{N} \left\{ \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} \frac{\partial f}{\partial b} + \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} \frac{\partial f}{\partial a} + \frac{\partial f}{\partial \mathbf{X}_i} \frac{\partial^2 f}{\partial a \partial b} + \frac{\partial^3 f}{\partial a \partial b \partial \mathbf{X}_i} f(a \mathbf{X}_j + b) - \dots \right. \\ \left. \frac{\partial E_\mu}{\partial a} \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} - \frac{\partial E_\mu}{\partial b} \frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} - \frac{\partial^2 E_\mu}{\partial a \partial b} \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^3 f}{\partial a \partial b \partial \mathbf{X}_i} \right\} \end{aligned} \quad (96)$$

$$\begin{aligned} \frac{\partial^3 E_\nu}{\partial b^2 \partial \mathbf{X}_i} = \frac{2}{N} \left\{ 2 \frac{\partial f}{\partial b} \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} + \frac{\partial f}{\partial \mathbf{X}_i} \frac{\partial^2 f}{\partial b^2} + f(a \mathbf{X}_j + b) \frac{\partial^3 f}{\partial b^2 \partial \mathbf{X}_i} - \dots \right. \\ \left. 2 \frac{\partial E_\mu}{\partial b} \frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} - \frac{\partial^2 E_\mu}{\partial b^2} \frac{\partial f}{\partial \mathbf{X}_i} - \left(\frac{1}{N} \sum_{j=1}^N f(a \mathbf{X}_j + b) \right) \frac{\partial^3 f}{\partial b^2 \partial \mathbf{X}_i} \right\}. \end{aligned} \quad (97)$$

4.1.3 Gradient of CAVE Base Functions w.r.t. \mathbf{X}_i

Lastly, we can find the additional function derivatives as

$$\frac{\partial f}{\partial \mathbf{X}_i} = \frac{\partial}{\partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f'(a \mathbf{X}_i + b) \cdot a \quad (98)$$

$$\frac{\partial^2 f}{\partial a \partial \mathbf{X}_i} = \frac{\partial^2}{\partial a \partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f''(a \mathbf{X}_i + b) \cdot a \mathbf{X}_i + f'(a \mathbf{X}_i + b) \quad (99)$$

$$\frac{\partial^2 f}{\partial b \partial \mathbf{X}_i} = \frac{\partial^2}{\partial b \partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f''(a \mathbf{X}_i + b) \cdot a \quad (100)$$

$$\frac{\partial^3 f}{\partial a^2 \partial \mathbf{X}_i} = \frac{\partial^3}{\partial a^2 \partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f'''(a \mathbf{X}_i + b) \cdot a \mathbf{X}_i^2 + f''(a \mathbf{X}_i + b) \cdot 2 \mathbf{X}_i \quad (101)$$

$$\frac{\partial^3 f}{\partial a \partial b \partial \mathbf{X}_i} = \frac{\partial^3}{\partial a \partial b \partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f'''(a \mathbf{X}_i + b) \cdot a \mathbf{X}_i + f''(a \mathbf{X}_i + b) \quad (102)$$

$$\frac{\partial^3 f}{\partial b^2 \partial \mathbf{X}_i} = \frac{\partial^3}{\partial b^2 \partial \mathbf{X}_i} [f(a \mathbf{X}_i + b)] = f'''(a \mathbf{X}_i + b) \cdot a. \quad (103)$$

All of these equations constitute the gradient of Newton's method.

4.2 Remaining Gradients

We derived the formula for the gradient of Newton's method for two variables, but there are still five other cases. We finish here with Newton's method with one variable, and gradient descent with both one and two variables.

Gradient descent only requires us to know some combination of $\frac{\partial^2 \mathcal{L}}{\partial a \partial \mathbf{X}_i}$, $\frac{\partial^2 \mathcal{L}}{\partial b \partial \mathbf{X}_i}$, $\frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial \mathbf{X}_i}$, and $\frac{\partial^2 \mathcal{L}_\mu}{\partial b \partial \mathbf{X}_i}$, all of which were found in the previous section, so we're done with gradient descent. Newton's method requires a little more effort. Recall that Newton's update for variance optimization is

$$\left(\frac{\partial \mathcal{L}_\nu}{\partial a} \right) \left(\frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} \right)^{-1}. \quad (104)$$

Using the quotient rule, the gradient of Newton's method for variance optimization is

$$\left(\frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} \frac{\partial^2 \mathcal{L}_\nu}{\partial a \partial \mathbf{X}_i} - \frac{\partial \mathcal{L}_\nu}{\partial a} \frac{\partial^3 \mathcal{L}_\nu}{\partial a^2 \partial \mathbf{X}_i} \right) \left(\frac{\partial^2 \mathcal{L}_\nu}{\partial a^2} \right)^{-2}. \quad (105)$$

Similarly for the mean, we have the optimization step as

$$\left(\frac{\partial \mathcal{L}_\mu}{\partial b} \right) \left(\frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} \right)^{-1} \quad (106)$$

and the gradient as

$$\left(\frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} \frac{\partial^2 \mathcal{L}_\mu}{\partial b \partial \mathbf{X}_i} - \frac{\partial \mathcal{L}_\mu}{\partial b} \frac{\partial^3 \mathcal{L}_\mu}{\partial b^2 \partial \mathbf{X}_i} \right) \left(\frac{\partial^2 \mathcal{L}_\mu}{\partial b^2} \right)^{-2}. \quad (107)$$