

This is the problem I'm trying to solve:

$$f(\mathbf{z}) = \frac{a}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + \frac{b}{2} \|\mathbf{P}\mathbf{z} - \mathbf{P}\mathbf{s}\|_2^2 \quad (1)$$

$$\arg \min_{\mathbf{z}} f(\mathbf{z}) \quad (2)$$

where \mathbf{P} is a nonlinear operator, specifically a quantized projection operator. That is, $\mathbf{P}(\cdot) = \mathbf{W}^T \text{Quantize}(\mathbf{W}\cdot)$, where $\mathbf{W}^T \mathbf{W}$ is a projection operator.

$$\text{Quantize}(\mathbf{y}) = \text{round}\left(\frac{\mathbf{y}}{\mathbf{q}}\right) * \mathbf{q} \quad (3)$$

(Division here is element-by-element).

In an early attempt to solve my problem, I pretended \mathbf{P} was a linear projection operator and came up with an approximate solution to my original problem.

$$\mathbf{z}_{\text{approx}} = \mathbf{x} + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) \quad (4)$$

Now, plugging $\mathbf{z} = \mathbf{z}_{\text{approx}} + \Delta\mathbf{z}$ back into the original equation, I have the expression:

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| \mathbf{P}\left(\mathbf{x} + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z}\right) - \mathbf{P}\mathbf{s} \right\|_2^2 \quad (5)$$

In analyzing this expression, I have found it helpful to define a function $\epsilon(\Delta\mathbf{z})$:

$$\epsilon(\Delta\mathbf{z}) = \mathbf{P}\left(\mathbf{x} + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z}\right) - \mathbf{P}\mathbf{x} - \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) \quad (6)$$

$$\mathbf{P}\left(\mathbf{x} + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z}\right) = \mathbf{P}\mathbf{x} - \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \epsilon(\Delta\mathbf{z}) \quad (7)$$

Returning to the objective function:

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| \mathbf{P}\mathbf{x} - \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \epsilon(\Delta\mathbf{z}) - \mathbf{P}\mathbf{s} \right\|_2^2 \quad (8)$$

Simplifying

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| -\frac{a+b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \epsilon(\Delta\mathbf{z}) \right\|_2^2 \quad (9)$$

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| -\frac{a}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \epsilon(\Delta\mathbf{z}) \right\|_2^2 \quad (10)$$

I have 2 objective terms, and it will help to give them names:

$$f_1(\Delta \mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b} (\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \Delta \mathbf{z} \right\|_2^2 \quad (11)$$

$$f_2(\Delta \mathbf{z}) = \frac{b}{2} \left\| -\frac{a}{a+b} (\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \epsilon(\Delta \mathbf{z}) \right\|_2^2 \quad (12)$$

$$f(\mathbf{z}) = f_1(\mathbf{z} - \mathbf{z}_{\text{approx}}) + f_2(\mathbf{z} - \mathbf{z}_{\text{approx}}) \quad (13)$$

Now for some observations:

1. Adding a component to $\Delta \mathbf{z}$ that is orthogonal to the span of the columns of \mathbf{W}^T increases the first term of the objective f_1 without affecting the second term f_2 .

$$f_1(\Delta \mathbf{z}) \geq f_1(\mathbf{W}^T \mathbf{W} \Delta \mathbf{z}) \quad (14)$$

$$f_2(\Delta \mathbf{z}) = f_2(\mathbf{W}^T \mathbf{W} \Delta \mathbf{z}) \quad (15)$$

Therefore,

$$(\mathbf{I} - \mathbf{W}^T \mathbf{W})(\Delta \mathbf{z})_{\text{optimal}} = 0 \quad (16)$$

2. In the simplified case of no quantization $\mathbf{P} = \mathbf{W}^T \mathbf{W}$:

$$\epsilon(\Delta \mathbf{z}) = \mathbf{W}^T \mathbf{W} \Delta \mathbf{z} \quad (17)$$

And so, if the quantization process is removed,

$$(\Delta \mathbf{z})_{\text{optimal}} = 0 \quad (18)$$

3. For $\alpha \in [0, 1]$:

$$f_2(\alpha \epsilon(0)) = f_2(0) \quad (19)$$

Furthermore, this is true even if I scale elements of $\mathbf{W}\epsilon(0)$ individually:

For $\alpha_i \in [0, 1]$:

$$f_2(\mathbf{W}^T \text{diag}(\boldsymbol{\alpha}) \mathbf{W} \epsilon(0)) = f_2(0) \quad (20)$$

4. There are certain choices for $\boldsymbol{\alpha}$ from the previous observation that will decrease the first objective term f_1 :

$$\alpha_i = \begin{cases} 1 & \text{sign}(\mathbf{e}_i^T \mathbf{W} \epsilon(0)) = -\text{sign}(\mathbf{e}_i^T \mathbf{W} (\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x})) \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

Using the $\boldsymbol{\alpha}$ defined above:

$$f_1(\mathbf{W}^T \text{diag}(\boldsymbol{\alpha}) \mathbf{W} \epsilon(0)) \leq f_1(0) \quad (22)$$

5. The last couple of operations have focused on decreasing f_1 without affecting f_2 . Here, I observe it is also possible to select a $\Delta \mathbf{z}$ that decreases f_2 by more than it decreases f_1 .

$$\beta_i = \begin{cases} 1 + \nu & \text{sign}(\mathbf{e}_i^T \mathbf{W} \boldsymbol{\epsilon}(\mathbf{q}/2)) = \text{sign}(\mathbf{e}_i^T \mathbf{W}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x})) \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (23)$$

where ν is an arbitrarily small number to ensure the rounding occurs in the proper direction.

Using the β defined above:

$$f_1(\mathbf{W}^T \text{diag}(\beta) \mathbf{W} \boldsymbol{\epsilon}(\frac{\mathbf{W}^T \mathbf{q}}{2})) - f_1(0) \leq f_2(0) - f_2(\mathbf{W}^T \text{diag}(\beta) \mathbf{W} \boldsymbol{\epsilon}(\frac{\mathbf{W}^T \mathbf{q}}{2})) \quad (24)$$

6. Finally, the last couple observations can be combined for the optimal solution:

$$(\Delta \mathbf{z})_{\text{optimal}} = \mathbf{W}^T \text{diag}(\beta) \mathbf{W} \boldsymbol{\epsilon}(\frac{\mathbf{W}^T \mathbf{q}}{2}) + \mathbf{W}^T \text{diag}(\alpha) \mathbf{W} \boldsymbol{\epsilon}(0) \quad (25)$$

where

$$\alpha_i = \begin{cases} 1 & \text{sign}(\mathbf{e}_i^T \mathbf{W} \boldsymbol{\epsilon}(0)) = -\text{sign}(\mathbf{e}_i^T \mathbf{W}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x})) \\ 0 & \text{otherwise} \end{cases} \quad (26)$$

$$\beta_i = \begin{cases} 1 + \nu & \text{sign}(\mathbf{e}_i^T \mathbf{W} \boldsymbol{\epsilon}(\mathbf{q}/2)) = \text{sign}(\mathbf{e}_i^T \mathbf{W}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x})) \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

$$\mathbf{z}_{\text{optimal}} = \mathbf{x} + \frac{b}{a+b}(\mathbf{P}\mathbf{s} - \mathbf{P}\mathbf{x}) + \mathbf{W}^T \text{diag}(\beta) \mathbf{W} \boldsymbol{\epsilon}(\frac{\mathbf{W}^T \mathbf{q}}{2}) + \mathbf{W}^T \text{diag}(\alpha) \mathbf{W} \boldsymbol{\epsilon}(0) \quad (28)$$

I still need to solve a slight variation of the above problem. Hopefully, the solution can be found in a similar way.

$$f(\mathbf{z}) = \frac{a}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + \frac{b}{2} \|\mathbf{P}\mathbf{z} + (1 - \mu)\mathbf{P}\mathbf{y} - (2 - \mu)\mathbf{P}\mathbf{s}\|_2^2 \quad (29)$$

$$\arg \min_{\mathbf{z}} f(\mathbf{z}) \quad (30)$$

$$\mathbf{z}_{\text{approx}} = \mathbf{x} + \frac{b}{a+b}((2 - \mu)\mathbf{P}\mathbf{s} - (1 - \mu)\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x}) \quad (31)$$

To prevent derivations from falling off the page:

$$\mathbf{r} = (2 - \mu)\mathbf{P}\mathbf{s} - (1 - \mu)\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x} \quad (32)$$

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| \mathbf{P} \left(\mathbf{x} + \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right) + (1-\mu)\mathbf{P}\mathbf{y} - (2-\mu)\mathbf{P}\mathbf{s} \right\|_2^2 \quad (33)$$

$$\epsilon(\Delta\mathbf{z}) = \mathbf{P} \left(\mathbf{x} + \frac{b}{a+b} ((2-\mu)\mathbf{P}\mathbf{s} - (1-\mu)\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x}) + \Delta\mathbf{z} \right) - \mathbf{P}\mathbf{x} - \frac{b}{a+b} ((2-\mu)\mathbf{P}\mathbf{s} - (1-\mu)\mathbf{P}\mathbf{y} - \mathbf{P}\mathbf{x}) \quad (34)$$

$$\epsilon(\Delta\mathbf{z}) = \mathbf{P} \left(\mathbf{x} + \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right) - \mathbf{P}\mathbf{x} - \frac{b}{a+b} \mathbf{r} \quad (35)$$

$$\mathbf{P} \left(\mathbf{x} + \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right) = \mathbf{P}\mathbf{x} + \frac{b}{a+b} \mathbf{r} + \epsilon(\Delta\mathbf{z}) \quad (36)$$

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| \mathbf{P}\mathbf{x} + \frac{b}{a+b} \mathbf{r} + \epsilon(\Delta\mathbf{z}) + (1-\mu)\mathbf{P}\mathbf{y} - (2-\mu)\mathbf{P}\mathbf{s} \right\|_2^2 \quad (37)$$

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| -\mathbf{r} + \frac{b}{a+b} \mathbf{r} + \epsilon(\Delta\mathbf{z}) \right\|_2^2 \quad (38)$$

$$f(\mathbf{z}_{\text{approx}} + \Delta\mathbf{z}) = \frac{a}{2} \left\| \frac{b}{a+b} \mathbf{r} + \Delta\mathbf{z} \right\|_2^2 + \frac{b}{2} \left\| -\frac{a}{a+b} \mathbf{r} + \epsilon(\Delta\mathbf{z}) \right\|_2^2 \quad (39)$$

The important distinction here is that μ prevents \mathbf{r} from being quantized by \mathbf{P} 's quantization. So, rounding in the "right direction" could go too far.