- **Protein Quantification**

Let $G_i = \{P_1, P_2, \cdots, P_k\}$ be the proteins in the proteome, $i, k \in \mathbb{N}^+$

Define $P = \{p_1, p_2, \cdots, p_n\}$ as the identified peptides. $n \in \mathbb{N}^+$

Let $S_i$ be the quantitative intensity for each peptide $p_i$.

Note that $p_i$ is generated from $G_i$.

★ Classic Method of Weighted Mean

$S = \sum_{i=1}^{n} w_i S_i$  with  $\sum w_i = 1$  where $i \in \mathbb{N}^+$ and $w_i$ is the weight of the peptide.

↓

But it doesn't consider the uncertainty of observation !!!

↓

Common Sense:  Variance represents uncertainty

↓

Volatility

$$Var(x) = E[(X - \mu)^2]$$

↓

$Var(x) \uparrow \rightarrow$ Unreliable $\uparrow$

We can't get the true variance from Mass Spectrometry ō∆ō

because MS is not an ideal, infinite, repeatable measurement.

><
‿  You said you'd rely on "repetitions" to estimate the variance?

Nope!!! The estimation $= \sum($var of noise $+$ var of biological var $+$ instability of instrument$)$

↓

Quant UMS = Quantification Using Minimum Uncertainty in MS

[ On the premise of peptide-protein mapping, using all available signals
for the optimal protein quantification and make this quantity the most reliable. ]

↓ ?ⓞⓞ   Which one is reliable? Which one is not reliable?

⇩

Uncertainty → Variance Modeling

# Logic:

① Background: Do protein quantification for a precursor $p$.

Let $s_1, \ldots, s_N$ be the signal for either MS1 or MS2 fragmentation ion,
$N \in \mathbb{N}^+$.

Let $x_i \in \mathbb{R}^+$ be the observed intensity and let $C_i \in [0,1]$ be the score.

$x_i$ is $\log_2$ ratio.

Let $\delta$ represent Variance.

Goal: Construct the optimal estimation $\hat{x}$ to express the true value of precursor.

$$\hat{x} = \frac{\sum_{i=1}^{N} w_i x_i}{\sum_{i=1}^{N} w_i} \qquad \text{where} \quad w_i \propto \frac{1}{\delta_i^2} \qquad (1)$$

↓ How to build model for $\delta$?

② log-scale linear model:

$$\log(\delta_i^2) = \theta_1 Ts^{(i)} + \theta_2 Tc^{(i)} + \theta_3 Ts.c^{(i)}$$

where $Ts^{(i)} = S_i^{-1}$, $Tc^{(i)} = 1 - \sqrt{C_i}$, $Ts.c^{(i)} = \sqrt{Ts^{(i)} \cdot Tc^{(i)}}$, $\theta_k = (\theta_k')^2$, $k = 1,2,3$

Then $\delta_i^2 = e^{\theta_1 Ts^{(i)} + \theta_2 Tc^{(i)} + \theta_3 Ts.c^{(i)}}$

Thus $Wi = (\delta_i^2)^{-1} = e^{-(\theta_1 Ts^{(i)} + \theta_2 Tc^{(i)} + \theta_3 Ts.c^{(i)})} = f_i(\theta_1, \theta_2, \theta_3)$ (2)

③ Bias correction for both $X_i$ and $\hat{X}_i$

$$B_i(\alpha_1, \alpha_2, \alpha_3) = \alpha_1 \sqrt{Ts^{(i)}} + \alpha_2 \sqrt{Tc^{(i)}} + \alpha_3 \sqrt{Ts.c^{(i)}}$$

$$X_i' = X_i + B_i(\alpha_1, \alpha_2, \alpha_3)$$

④ Fusion

plug (2) into (1),

1. Direct Fusion Method

$$\hat{X}_i(\theta_1, \theta_2, \theta_3) = \frac{\sum_{i=1}^{N} Wi \, X_i}{\sum_{i=1}^{N} Wi}$$

$$= \frac{\sum_{i=1}^{N} f_i(\theta_1, \theta_2, \theta_3) \cdot X_i}{\sum_{i=1}^{N} f_i(\theta_1, \theta_2, \theta_3)}$$

$$= \frac{\sum_{i=1}^{N} e^{-(\theta_1 Ts^{(i)} + \theta_2 Tc^{(i)} + \theta_3 Ts.c^{(i)})} \cdot X_i}{\sum_{i=1}^{N} e^{-(\theta_1 Ts^{(i)} + \theta_2 Tc^{(i)} + \theta_3 Ts.c^{(i)})}} \qquad (3)$$

2. Iterative propagation across acquisition

For each fragmentation $m$, from run $j$ refer to run $j$

$$X_{m,j,i} = q_i + S_{m,j} - S_{m,i}$$

Weighted average for $i$: $X_{m,j} = \dfrac{\sum\limits_{i} W_{m,j,i} \cdot X_{m,j,i}}{\sum\limits_{i} W_{m,j,i}}$

Weighted average for $m$: $\bar{X}_j = \dfrac{\sum\limits_{m} W_m \cdot X_{m,j}}{\sum\limits_{m} W_m}$

⑤

1. Loss function <span style="color:red">classic</span>

$$\ell(\theta_1, \theta_2, \theta_3, \alpha_1, \alpha_2, \alpha_3) = \sum_{P} \sum_{i=1}^{N} f_i(\theta_1, \theta_2, \theta_3) \cdot \left[ x_i' - \hat{X}_i(\theta_1, \theta_2, \theta_3, \alpha_1, \alpha_2, \alpha_3) \right]^2$$

$$\ell(\theta, \alpha) = \sum_{P} \sum_{i=1}^{N} f_i(\theta) \cdot \left[ x_i + B_i(\alpha) - \hat{x}_i(\theta, \alpha) \right]^2$$

$\downarrow$

Weighted least squares

In conclusion, if $x_i$ and $\hat{x}_i$ have large difference, then $W_i$ is small

Essentially,     Self-supervised regression model + Weighted average model

2. Loss function - <span style="color:red">Quant UMS</span>

- Precision loss

$$\ell_{prec} = \frac{1}{N} \sum_{m,i} (X_{m,i} - S_{m_{sel},i})^2$$

where $X_{m,i}$ : The estimate value of feature m at aquisition i

$S_{m_{sel},i}$ : the signal of best m at run i

- Accuracy loss:

deviation        relative abundance

$$\ell_{acc} = \frac{\sum_{m,i} (X_{m,i} - S_{m_{sel},i})(S_{m_{sel},i} - \mu) W_i}{\sum_{m,i} W_i^2 \cdot Var} \quad ,$$

Where $S_i \sim N(\mu, Var)$,

$$W_i' = e^{-S_{m_{sel},i}}$$

- Class problem :   Signal with low abundance are weak

and easy to be ignore.

↓

Without any correction, weak signals has

been "averaged out".

↓

"Ratio Compression"

Ex:   low abunce $(S_{m_{sel},i} < \mu)$

↙                        ↓

$X_{m,i} - S_{m_{sel},i}$         $S_{m_{sel},i} - \mu$

$< 0$                         $< 0$

×

$\downarrow$

$>0$

High abunce $\rightarrow$ product $>0$

$\downarrow$

$W_i \uparrow$ when signal is weaker

$\downarrow$

- Optimize $L_{acc}$ :

Low signal will not be undervalued

High                    overvalued

$\rightarrow$ formula $\approx 0$

In conclusion,

$$\mathcal{L}(\theta, \alpha) = l_{prec} + \lambda \, l_{acc}, \quad \lambda \text{ is a constant}$$

$$\nabla \mathcal{L} = \nabla_\theta l_{prec} + \lambda \nabla_\theta l_{acc}$$

$$\nabla \mathcal{L} = \nabla_\alpha l_{prec} + \lambda \nabla \, l_{acc}$$

⑥ Backpropagation, update $\theta$, $a$

Goal: minimize $\mathcal{L}(\theta, \alpha)$

$$\frac{\partial \mathcal{L}}{\partial \theta_k} \qquad\qquad \frac{\partial \mathcal{L}}{\partial \alpha_k}$$

$$k = 1, 2, 3$$

- Gradient descent:

$$\theta_k^{(t+1)} = \theta_k^{(t)} - \eta \cdot \frac{\partial \mathcal{L}}{\partial \theta_k}$$

$$\alpha_k^{(t+1)} = \alpha_k^{(t)} - \eta \cdot \frac{\partial \mathcal{L}}{\partial \alpha_k}, \qquad \eta \text{ is the learning rate}$$

↓ I know I need to follow the current gradient, but how far I should go? ( How large is the $\eta$ ?)

↓

Automatically adjust $\eta$ !!!

↓

- Armijo Condition:

$$\mathcal{L}(x - \eta \nabla \mathcal{L}) \leq \mathcal{L}(x) - c \cdot \eta \cdot \|\nabla \mathcal{L}\|^2$$

where $\mathcal{L}(x)$ : Loss function

$\quad x$ : $\theta$ or $a$

$\quad d$ : $-1$

$\quad \eta$ : learning rate

$$\nabla f(x)^T : \frac{\partial L}{\partial \theta_k} \text{ or } \frac{\partial L}{\partial \alpha_k}$$

$C$ : constant — Control the tolerance level of declining

## Euclidean Norm:

$$\| \nabla f(\theta, \alpha) \|^2 = \sum_i \left( \frac{\partial f}{\partial \theta_i} \right)^2 + \sum_i \left( \frac{\partial f}{\partial \alpha_i} \right)^2$$

↓

## Why we use Armijo condition?

∵ $L(\theta, \alpha)$ is non-linear and has nested structure

∴ It's very easy that big $\eta$ cause the increase of loss.

∴ We need to calculate $\eta$.

Thus, input: $X_i$, $\theta$. $\alpha$

↓

linear model: $\delta^2$. $w$. bias

fusion: $\hat{x}$

Loss func: $L(\theta, \alpha)$

↓

backpropagation: $\frac{\partial L}{\partial \theta_k}$ , $\frac{\partial L}{\partial \alpha_k}$

Armijo : $\eta$

↓

update $\theta$. $\alpha$

↓

Until $L(\theta, \alpha)$ is minimized and converged.

## Explaination for Fig 1.

|   | Human | E.coli |
|---|-------|--------|
| A | 1     | 50     |
| B | 1     | 33     |
| C | 1     | 20     |

| A/C | Human | E.coli |
|-----|-------|--------|
| FC  | 1     | 2.5    |
| Fold change | | |
| $log_2 FC$ | 0 | 1.32 |