

author: Ocean Wong  
(Hoi Yeung Wong)  
supervisor: Chantal Nobs, Robin Smith  
organization: Culham Centre for Fusion Energy  
date: October 2019

---

## Abstract

In this document, the method of particle spectrum unfolding using Tikhonov Regularization is presented. Kullback-Leibler divergence was chosen as the regularizing function. The results has been implemented as a class within the neutron spectrum unfolding suite.

The same framework of minimizing the loss function, linearization of loss function gradient landscape, and proving the uniqueness (or non-uniqueness) of the global minimum, can be re-used by other researches who wishes to create a different regularization program based on another (or multiple others) metric(s).

A general neutron spectrum unfolding algorithm by regularization using relative-entropy

## 1 Background

regularization algorithm uses a metric, which is a linear mixture of  $(\chi^2)_{\text{reactionrates}}$  and the deviation of the solution spectrum from the *a priori* spectrum. The mixing ratio between these two quantities is controlled by the Tikhonov parameter  $\tau$ , which more or less weight on the relative-entropy.

This method of particle spectra unfolding is applicable to any general unfolding problem, as nothing except the reaction rates, response matrix, and *a priori* is known. This is particularly useful for fusion neutron spectrum unfolding, where there are not many experimental data available, thus no reasonable parametrisation or other sophisticated method of simplifying the neutron spectrum can be used.

$$\text{loss}(\phi_{\text{sol}}) = (\chi^2)_{\text{reaction rates}}(\phi_{\text{sol}}) + \tau \cdot \text{RelativeEntropy}_{\phi_0}(\phi_{\text{sol}}) \quad (1)$$

where

$\phi_0$  is the *a priori* spectrum,

$\phi_{\text{sol}}$  is the algorithm's output spectrum,

$\text{Loss}(\phi_{\text{sol}})$  is the value outputted by the loss function, sometimes called the loss function

Due to the underdetermination condition of the unfolding problem, there exist multiple (i.e. infinitely many) solutions which will minimize  $\chi^2$  (potentially down to zero). By also requiring it to deviate minimally from the *a priori* spectrum, a unique solution to the problem can be found instead.

### 1.1 Motivation

While some programs of particle spectra unfolding via regularization exist, it is not rigorously derived and do not allow the user to specify a regularization constant. MAXED in particular, forces  $(\chi^2) = \text{number of degree of freedoms}$  [7]. This form of "regularization" leads to no insight to the user of how far off is the *a priori* spectra, and therefore if the spectrum is accurate or not. It also does not guarantee that

the solution is unique, since multiple (infinitely many) solutions with  $\chi^2 = 1$  can be found<sup>1</sup>; while for this program, only finitely many (only one) of such minima will be found. (See Appendix B) By carrying out the following derivation, which does not rely on any pre-existing numerical solvers to obtain a solution (i.e. no bin-hopping or simulated annealing), a minimum loss value in equation 1 is guaranteed, without the chance of being trapped at sub-optimal stationary points. A program is created using the result of this derivation.

## 2 Justification of loss function (relative-entropy + $(\chi^2)$ ) chosen

$\chi^2$  was used to quantify the deviation of the unfolded solution neutron spectrum's would-have-been reaction rate from the measured reaction rates. This is a widely used metric for measuring the deviation between pairs of numbers that are otherwise uncorrelated, and is generally accepted as the standard metric used to measure the deviation between reaction rates [4][8].

Perhaps a more interesting question is why should we use relative-entropy as the regularizing function in the loss function.

The regularizing function, in the context of neutron spectrum unfolding, refers to the function used to compute the “deviation” of one spectrum from another (from the solution spectrum to the *a priori* spectrum).

When choosing a regularizing function, it is desirable to have the following features:

1. Can naturally obey the physical requirement of having no negative flux without applying additional external constraint onto the set of equations;
2. For the same amount of absolute deviation  $\Delta\phi$ , harsher penalty should be given to bin with lower *a priori* flux. Penalty contributed by the  $i$ -th bin; penalty contributed by the  $j^{th}$  bin if  $|\phi_{sol,i} - \phi_{0,i}| = \Delta\phi = |\phi_{sol,j} - \phi_{0,j}|$

To satisfy both of these requirements,

1. The gradient of the loss function must approach  $-\infty$  as the flux approaches 0,
2. The magnitude of the 2<sup>nd</sup> order derivative of the loss function  $\frac{\partial^2(\text{loss})}{\partial\phi_{sol,i}^2}$  must increase with decreasing *a priori* value  $\phi_{0,i}$ <sup>2</sup>

Other metrics considered before writing this document includes:

- Fractional deviation (obeys requirement 2 but not 1)
- Mean squared deviation in log-space
- Absolute deviation in log-space

---

<sup>1</sup>MAXED deals with this problem of infinitely-many-possible-solutions by opening up an extra degree of freedom, removing the  $\sum \phi_i = 1$  constraint, seeking to minimize the deviation of this sum  $\sum \phi_i$  from 1 instead. But, as the author of this paper, I disagree with that approach, as it invalidates the premise of having  $f_i$  as a “normalized probability distribution”.

<sup>2</sup>If the reader disagree with requirement 2 but still would like a method that fulfill requirement 1, then s/he can use a regularizing function with two Tikhonov parameters: Regularizing function =  $\tau_1(\chi^2)_{\phi,\phi_0} + \tau_2 S(\phi)$ ; where  $S(\phi)$  is the self-entropy of the probability distribution of the solution spectrum  $\phi$ , and  $(\chi^2)_{\phi,\phi_0}$  is the usual definition of  $\chi^2$  (sum of square of weighted deviations of flux in each bin, between the solution spectrum and the *a priori* spectrum.)

- Determinant of the Fisher Information Matrix <sup>3</sup>

Relative entropy satisfies both of the requirements above, with the additional benefit of being the only unique and non-self-contradictory metric that can be used to measure deviation between probability distributions (See Appendix A in [6]). Therefore it is chosen as the metric to be used in the regularizing function in this paper.

Serendipitously, the loss function, when composed in this manner, has no local trapping suboptimal points (See appendix B), allowing to use of Lagrangian multipliers to find the point of minimum minimization, which is detailed in section 4.

There are other implementations of the regularization method in ways that does not require an *a priori*. Take minimum (self-)entropy for example, which is the algorithm used by [1] and [3]. This is only a special implementation of the principle of MaxEnt (Maximum (relative-)Entropy), where a naïve *a priori* is assumed, i.e. it is equivalent to using a flat *a priori* in MAXED.

It is also worth noting that the commonly used neutron spectrum unfolding code GRAVEL obeys both requirements, even though it is not program that uses a regularization code. It takes steps in the log-flux space (thus fulfilling requirement ??) in order to minimize the  $\chi^2$  of the reaction rates; and the starting point in the log-flux space is the *a priori* (thus fulfilling requirement 2).

## 2.1 Definition of relative-entropy

Before proceeding any further, it is useful to know what relative-entropy is. It is also known as Kullback-Leibler divergence ( $D_{KL}$ ); and has been misnamed as cross-entropy( $H$ ) in some papers[6] (See equation 5 for details). Both measures the deviation of one normalized probability distribution from another. These probability distribution can be continuous or discrete.

For example, let there be two discrete probability distributions Q and P , over x:

$$H(P, Q) = - \sum_x P(x) \log(Q(x)) \quad (2)$$

$$S(P) = - \sum_x P(x) \log(P(x)) \quad (3)$$

$$D_{KL}(P||Q) = \sum_x P(x) \log\left(\frac{P(x)}{Q(x)}\right) \quad (4)$$

$$D_{KL}(P||Q) = H(P, Q) - S(P) \quad (5)$$

where  $S(P)$  is the (self-)entropy of the distribution. (Intuitively, when analogized to drawing balls from a bag,  $S(P)$  is how “surprising” the resulting ball’s colour is on average.) When  $P$  is fixed (as is the *a priori* ), then  $S(P)$  is a constant, so  $D_{KL}(P||Q)$  and  $H(P, Q)$  only differ by a fixed, constant value for any P.

Note the non-commutativity of  $D_{KL}(P||Q)$ : if the role of P and Q in equation 4 were reversed, a different answer will be obtained; but it will still give a loss function landscape that only has one global minimum.

The choice of the base of the  $\log$  in the equations above is arbitrary. (In the ball-drawing analogy this number is equal to the number of colours available.) In

---

<sup>3</sup>Fisher information of a variable vector (which the neutron spectrum is) should form a MATRIX, not a scalar, unlike the papers [2][5] mentioned.

the context of information theory and computer science this is usually chosen as 2, and given the unit “shannons”; but for physics application, and for the rest of this text, it will be chosen as base e, and given the unit “nats”, where 1 shannon = 0.693 nats.

### 3 Assumption

Let’s assume that we have a response matrix  $\underline{\mathbf{R}}$  where each element is known to sufficiently high precision.  $R_{ki}$  refers to the probability of a nuclide of type  $k$  being created upon the incidence of a neutron with energy within energy bin  $i$ .

#### 3.1 Notations

In previous texts, the neutron flux is usually represented with  $\phi$ , and the reaction rates with  $Z$ . In this text, however, the neutron flux needs to be normalized to obtain a probability distribution instead. The probability distribution of neutron will be indicated with  $\mathbf{f}$ ,

$$\sum_i f_i = 1 \quad (6)$$

in line with the notations used in [6] derivation. Additionally, the reaction rate will not be use; instead the following quantity  $\mathbf{N}$

$$N_k = \# \text{ of nuclide } k \text{ created per incident neutron} \quad (7)$$

$$= \frac{\text{total \# of nuclide } k \text{ created}}{\text{total n. yield}} = \frac{\text{reaction rate } k}{\text{n. per second}} = \frac{Z_k}{\text{n. per second}} \quad (8)$$

will be use (contrary to previous texts which may use  $N$  as the symbol or no. of nuclides in total).

To make the following derivation easier, we assume that the measurements of “# of nuclide  $k$  created per incident neutron”  $N_k$  are mutually independent, forming a vector  $\mathbf{N}$  of with a covariance matrix which is purely diagonal,

$$\underline{\underline{\mathbf{S}_{\mathbf{N}}}} = \begin{pmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_k & & \\ & & & \ddots & \\ & & & & \sigma_m \end{pmatrix} \quad (9)$$

However, if one wishes to account for the covariances between different measured  $N_k$ , they can repeat the following derivation with  $(\chi^2) = (\mathbf{N}' - \mathbf{N}) \cdot \underline{\underline{\mathbf{S}_{\mathbf{N}}}} (\mathbf{N}' - \mathbf{N})$

instead of  $(\chi^2) = \sum_k (\frac{N'_k - N_k}{\sigma_k})^2$ . One of the possible reasons to do so is to account for the covariance in uncertainty between  $N_k$  due to uncertainty of the total n. yield, which is the common denominator to all of these reaction rates. Even if  $Z_k \{ \forall 1 \leq k \leq m \}$  are all obtained via independent measurements, they need to be divided by the same neutron yield as shown in equation 8, introducing a large component of positive covariance between them.

Below are some of the notations which will be used:

number of reactions =	$m$
	so that daughter nuclide $k$ has a range $1 \leq k \leq m$
number of energy groups =	$i$
	so that the $i^{th}$ or $j^{th}$ energy group $1 \leq i, j \leq n$
<u>A priori</u> probability distribution of neutron spectrum =	$\mathbf{f}^{DEF}$
probability distribution of neutron spectrum =	$\mathbf{f}$ (vectors are written in bold italics)
Measured # of nuclide created per incident neutron =	$\mathbf{N}$ (without prime ('))
Response matrix =	$\underline{\underline{\mathbf{R}}}$
element of a matrix is indicated with a pair of subscript indices, e.g. =	$R_{ij}$
element of a vector is indicated with a subscript index, e.g. =	$f_i$
Iteration of a vector/matrix is indicated with superscript, e.g. =	$\mathbf{v}^g$

The definitions of some more specific terms will be made clear further down the derivation. See equation 33, 34, and 66 for the definition of  $\mathbf{c}$ ,  $\underline{\underline{\mathbf{M}}}$ , and  $\underline{\underline{\mathbf{J}}}$ .

## 4 Derivation

### 4.1 Formulation (Lagrangian Multiplier)

For a fixed value of  $\tau$  the following m+1 constraints has to be followed:

$$N'_k = \sum_i^n R_{ki} f_i \quad \text{for } 1 \leq k \leq m \quad (10)$$

$$\sum_i f_i = 1 \quad (11)$$

where  $N'_k$  is the hypothetical number of nuclide  $k$  generated per incident neutron given the hypothetical neutron spectrum spectrum  $\mathbf{f}$ .

And the following quantities has to be minimized:

$$\tau(D_{KL}(\mathbf{f}, \mathbf{f}^{DEF})) = \tau \sum_i f_i \ln\left(\frac{f_i}{f_i^{DEF}}\right) \quad (12)$$

$$\chi^2(\mathbf{N}', \mathbf{N}) = \sum_k \left(\frac{N'_k - N_k}{\sigma_k}\right)^2 \quad (13)$$

Note that  $D_{KL}$  is always positive when constraint 11 is obeyed; and reaches a minimum when  $\mathbf{f} = \mathbf{f}^{DEF}$ .

Thus the Lagrangian can be constructed as

$$\mathcal{L}(\mathbf{f}, \boldsymbol{\lambda}, \mathbf{N}', \mu) = \tau \sum_i f_i \ln\left(\frac{f_i}{f_i^{DEF}}\right) + \sum_k \left(\frac{N'_k - N_k}{\sigma_k}\right)^2 + \sum_k \lambda_k [N'_k - \sum_i^n R_{ki} f_i] + \mu [\sum_i f_i - 1] \quad (14)$$

where  $\boldsymbol{\lambda}$  is a list of Lagrangian multipliers  $\lambda_k$ ,  $1 \leq k \leq m$ ; and  $\mu$  is a Lagrangian multiplier.

To obtain the minimum loss value where  $loss = \chi^2 + \tau(D_{KL})$ ,

$$\nabla \mathcal{L} = 0 \quad (15)$$

i.e.

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial f_i} = 0 & \text{for } 1 \leq i \leq n \end{cases} \quad (16)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial N'_k} = 0 & \text{for } 1 \leq k \leq m \end{cases} \quad (17)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \lambda_k} = 0 & \text{for } 1 \leq k \leq m \end{cases} \quad (18)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \mu} = 0 \end{cases} \quad (19)$$

This leads to the following  $n+2m+1$  simultaneous equations:

$$\begin{cases} \tau \left[ \ln\left(\frac{f_i}{f_i^{DEF}}\right) + 1 \right] + \mu + \sum_k^m \lambda_k R_{ki} = 0 & \text{for } 1 \leq i \leq n \end{cases} \quad (20)$$

$$\begin{cases} 2 \frac{N'_k - N_k}{\sigma_k^2} - \lambda_k = 0 & \text{for } 1 \leq k \leq m \end{cases} \quad (21)$$

$$\begin{cases} \sum_j^n R_{kj} f_j = N'_k & \text{for } 1 \leq k \leq m \end{cases} \quad (22)$$

$$\begin{cases} \sum_i f_i = 1 \end{cases} \quad (23)$$

Substituting the 3<sup>rd</sup> line into the 2<sup>nd</sup> line removes  $N'_k$  from this set of equations entirely; then substituting the 2<sup>nd</sup> line into the 1<sup>st</sup> removes the  $n$  Lagrangian multipliers  $\lambda_k$  from this set of equations entirely, leaving only a set of equations dependent on the variables  $f_i$  ( $1 \leq i \leq n$ ) and  $\mu$ , as well as  $n+1$  linearly-independent equations.

$$\begin{cases} \tau \left[ \ln\left(\frac{f_i}{f_i^{DEF}}\right) + 1 \right] + \mu + \sum_k^m \left[ 2 \frac{R_{ki}}{\sigma_k^2} \left( \left( \sum_j^n R_{kj} f_j \right) - N_k \right) \right] = 0 \end{cases} \quad (24)$$

$$\begin{cases} \sum_j^n f_j = 1 \end{cases} \quad (25)$$

Rearrange the top line (note: each  $i$  represent 1 equation,  $1 \leq i \leq n$ , so the top line entails  $n$  equations), such that all linear terms of  $f_i$  and  $\mu$  stays on the LHS, while the constants and non-linear terms in  $f_i$  moves to the RHS.

$$\sum_k^m \left[ 2 \frac{R_{ki}}{\sigma_k^2} \left( \sum_j^n R_{kj} f_j \right) \right] - \sum_k^m \left( 2 \frac{N_k R_{ki}}{\sigma_k^2} \right) + \mu = -\tau \left[ \ln \left( \frac{f_i}{f_i^{DEF}} \right) + 1 \right] \quad (26)$$

$$\sum_k^m \left[ 2 \frac{R_{ki}}{\sigma_k^2} \left( \sum_j^n R_{kj} f_j \right) \right] + \mu = -\tau \left[ \ln \left( \frac{f_i}{f_i^{DEF}} \right) + 1 \right] + \sum_k^m \left( 2 \frac{N_k R_{ki}}{\sigma_k^2} \right) \quad (27)$$

Transposing the order of summation in the first term yields

$$\sum_k^m 2 \frac{R_{ki}}{\sigma_k^2} \sum_j^n (R_{kj} f_j) = \sum_k^m \sum_j^n \left[ \left( 2 \frac{R_{ki} R_{kj}}{\sigma_k^2} \right) f_j \right] = \sum_j^n \left[ \left( 2 \sum_k^m \left( \frac{R_{ki} R_{kj}}{\sigma_k^2} \right) \right) f_j \right] \quad (28)$$

Now if we construct matrix  $\underline{\mathbf{A}}$  such that for  $1 \leq i, j \leq n$

$$A_{ij} = 2 \sum_k^m \left( \frac{R_{ki} R_{kj}}{\sigma_k^2} \right) \quad (29)$$

And further rearranging so only the constants are on the RHS, then the set of simultaneous equations can be expressed as:

$$\left\{ \begin{array}{l} \sum_j^n A_{ij} f_j + \tau \ln(f_i) + \mu = \tau [\ln(f_i^{DEF}) - 1] + 2 \sum_k^m \left( \frac{N_k R_{ki}}{\sigma_k^2} \right) \\ \sum_j^n f_j = 1 \end{array} \right. \quad (30)$$

$$\sum_j^n f_j = 1 \quad (31)$$

## 4.2 How not to solve it

One may attempt to decouple the dependence of the  $i^{th}$  term from itself in equations 30, perhaps by moving all terms dependent on  $f_i$  in the  $i^{th}$  equation to the LHS, and all terms dependent on  $f_{j \neq i}$  to the RHS, so that each line can be expressed as  $F_{i/LHS}(f_i) = F_{i/RHS}(f_{j \neq i}, constants)$ , i.e. equating the result of non-linear function  $F_{i/LHS}(f_i)$  to a the result of a (linear) function dependent on  $f_j$ .

This seems promising at first, until one tries to solve for  $f_1$ , and realising that s/he still knows nothing about the values on the RHS as  $f_j$  (for  $j=2, \dots, n$ ) are unknown.  $f_2$  is a function of  $f'_j$  where  $j' = 1, 3, 4, \dots, n$ , etc. To evaluate  $f_2$  will lead to circular reference towards  $f_1$ , leading to infinite recursion.

It is also tempting to give up at this stage and allow a generic minimization algorithm (e.g. L-BFGS) to find the solution of these equations. However, due to the logarithmic term present, most algorithm will inevitably overshoot into the  $f_i < 0$  region when exploring the solution space, leading to failure.

## 4.3 Linearization

Instead, a much more elegant technique is stated as follows.

In general, in a set of fully-determined simultaneous equations, the variables required can be stacked into a column vector  $\mathbf{v}$ , and the constant terms can be moved to the RHS to be stacked into a column vector  $\mathbf{c}$ . Assuming that only linear operations on  $\mathbf{v}$  is present in the system, then these simultaneous equations can be

expressed in linear algebra notations as  $\underline{\underline{\mathbf{M}}}\mathbf{v} = \mathbf{c}$ , where  $\underline{\underline{\mathbf{M}}}$  is a non-singular square matrix, so  $\mathbf{v} = \underline{\underline{\mathbf{M}}}^{-1}\mathbf{c}$ .

For the regularization problem described by equation 30 and 31, the set of simultaneous equations is fully determined, but not expressible as a linear operator  $\underline{\underline{\mathbf{M}}}$  operating on a vector, as it involves the non-linear operator  $\ln$ .

One way of getting around this is via Taylor expansion. By incorporating the first term of the Taylor expansion into the LHS of equation 30 and ignoring all higher order terms, an approximate solution can be obtained.

This process can be repeatedly applied onto the result obtained to iteratively improve the accuracy of this approximation until it is much higher than the precision of the floating point.

To more succinctly express the set of equations above (30 and 31), we will relabel them into a matrix multiplication equation,  $\underline{\underline{\mathbf{M}}}\mathbf{v} = \mathbf{c}$ .

The most obvious next step is to turn the LHS of 30 and 31  $\underline{\underline{\mathbf{M}}}\mathbf{v}$  where  $\underline{\underline{\mathbf{M}}}$  has size  $(n+1) \times (n+1)$ , variable vector  $\mathbf{v}$  has length  $(n+1)$ , and the RHS into a constant vector  $\mathbf{c}$  has length  $(n+1)$ . However, we can eliminate the variable  $\mu$  from the set of equations by subtracting the  $(i+1)^{th}$  line from the  $i^{th}$  line. Note that  $\mu$  has no immediately obvious physical meaning (it represents the steepness of the function in the direction normal to the constraint hyperplane), so it is not of interest to us. Then the set of  $n$  simultaneous equations in (30) can be reduced to the following  $n-1$  equations, where  $1 \leq i \leq n-1$ :

$$\sum_j^n [A_{ij} - A_{(i+1)j}] f_j + \tau [\ln(f_i) - \ln(f_{i+1})] = \tau [\ln(f_i^{DEF}) - \ln(f_{i+1}^{DEF})] + 2 \sum_k^m \left[ \frac{N_k(R_{ki} - R_{k(i+1)})}{\sigma_k^2} \right] \quad (32)$$

Such that equation 31 and 32 together forms  $n$  simultaneous equations in  $n$  variables, i.e. a fully determined system.

In attempt to condense them further into vector and matrix representations, we can use the following notations:

$$M_{ij} = \begin{cases} A_{ij} - A_{(i+1)j} & \text{for } 1 \leq i \leq n-1, 1 \leq j \leq n \\ 1 & \text{for } i = n, 1 \leq j \leq n \end{cases} \quad (33)$$

$$c_i = \begin{cases} \tau [\ln(f_i^{DEF}) - \ln(f_{i+1}^{DEF})] + 2 \sum_k^m \left( \frac{N_k(R_{ki} - R_{k(i+1)})}{\sigma_k^2} \right) & \text{for } 1 \leq i \leq n \\ 1 & \text{for } i = n \end{cases} \quad (34)$$

such that

$$\begin{cases} \sum_j^n M_{ij} f_j + \tau [\ln(f_i) - \ln(f_{i+1})] & -c_i = 0 & \text{for } 1 \leq i \leq n-1 & (35) \\ \sum_j^n M_{ij} f_j & -c_i = 0 & \text{for } i = n & (36) \end{cases}$$

Express the  $i^{th}$  line of LHS of equations 35 and 36 as a function  $F_i$ . Define  $\mathbf{u}$  to be the test solution vector, while the  $\mathbf{v}$  is the correct solution vector, such that

$$F_i(\mathbf{u}) = 0 \quad (37)$$



which can be written in vector form

$$\mathbf{F}(\mathbf{u}) = \mathbf{0} \quad (38)$$

if  $\mathbf{u} = \text{solution} \equiv \mathbf{v}$ .

For convenience the vector  $\mathbf{\Delta}$  is defined:

$$\mathbf{\Delta} \equiv (\mathbf{v} - \mathbf{u}) \quad (39)$$

When searching for the solution, the true value of  $\mathbf{v}$  is not known; but guess values of  $\mathbf{u}$  can be used.

Consider the Taylor expansion of  $\ln(v_i)$  at  $\ln(u_i)$ ,

$$\ln(v_i) = \ln(u_i) + \frac{v_i - u_i}{u_i} - O(u_i, \Delta_i) \quad (40)$$

$$\tau \ln(v_i) = \frac{\tau}{u_i} v_i - \tau(-\ln(u_i) + 1 + O(u_i, \Delta_i)) \quad (41)$$

where the higher order terms  $O(u_i, \Delta_i)$

$$-O(u_i, \Delta_i) = -\frac{1}{2} \left( \frac{\Delta_i}{u_i} \right)^2 + \frac{1}{3} \left( \frac{\Delta_i}{u_i} \right)^3 - \frac{1}{4} \left( \frac{\Delta_i}{u_i} \right)^4 + \frac{1}{5} \left( \frac{\Delta_i}{u_i} \right)^5 \quad (42)$$

$$= -\sum_{l=2}^{\infty} \frac{1}{l} \left( \frac{\Delta_i}{u_i} \right)^l \quad (43)$$

Therefore Taylor expanding the logarithms in 35 gives:

$$\begin{aligned} \sum_j^n M_{ij} v_j + \frac{\tau}{u_i} v_i - \frac{\tau}{u_{i+1}} v_{i+1} \\ = c_i - \tau[\ln(u_i) - \ln(u_{i+1})] + \tau[O(u_i, \Delta_i) - O(u_{i+1}, \Delta_{i+1})] \end{aligned} \quad (44)$$

expressing the equation 44 in terms of  $\mathbf{\Delta}$  and  $\mathbf{u}$  only,

$$\begin{aligned} \sum_j^n M_{ij} \Delta_j + \frac{\tau}{u_i} \Delta_i - \frac{\tau}{u_{i+1}} \Delta_{i+1} + \sum_j^n M_{ij} u_j + \tau \left[ \frac{u_i}{u_i} - \frac{u_{i+1}}{u_{i+1}} \right] \\ = c_i - \tau[\ln(u_i) - \ln(u_{i+1})] + \tau[O(u_i, \Delta_i) - O(u_{i+1}, \Delta_{i+1})] \end{aligned} \quad (45)$$

Manipulate equation 45, leaving only terms linearly dependent on  $\mathbf{\Delta}$  on the left,

$$\begin{aligned} \sum_j^n M_{ij} \Delta_j + \frac{\tau}{u_i} \Delta_i - \frac{\tau}{u_{i+1}} \Delta_{i+1} \\ = c_i - \sum_j^n M_{ij} u_j - \tau[\ln(u_i) - \ln(u_{i+1})] + \tau[O(u_i, \Delta_i) - O(u_{i+1}, \Delta_{i+1})] \end{aligned} \quad (46)$$

Meanwhile, the  $n^{th}$  line of the system of simultaneous equations (equation 36) gives

$$\sum_j^n M_{ij} u_j - c_i = 0 \quad (47)$$

$$\sum_j^n M_{ij} v_j - c_i = 0 \quad (48)$$

$$\sum_j^n M_{ij} \delta_j = 0 \quad (49)$$

where  $i = n$ . Note  $M_{nj} = 1 \forall 1 \leq j \leq n$  and  $c_n = 1$  as defined in 33 and 34; so all of the  $M_{ij}$  and  $c_i$  in the three equations above are unity.

#### 4.4 Convergence condition and underrelaxation

Rewrite equation 39 as

$$\mathbf{v} = \mathbf{u} + \mathbf{\Delta} \quad (50)$$

Where  $\mathbf{\Delta}$  is the deviation of  $\mathbf{u}$  from the true root  $\mathbf{v}$  in the domain space. It is difficult to accurately calculate  $\mathbf{\Delta}$ . However, if we take the approximation of it up to the first order as shown in equation 44 to 46, and call it  $\mathbf{\delta}$ , we can rewrite equation 50 as

$$\mathbf{f}^g = \mathbf{f}^{g-1} + \alpha \mathbf{\delta}^{g-1} \quad (51)$$

for the  $g^{th}$  iteration, where  $\mathbf{\delta}^{g-1}$  is the vector proportional to the difference between the solution vector (i.e. solution flux) at the (g-1)-th and g-th iterations.

Note that we have introduced an underrelaxation factor  $0 < \alpha \leq 1$  to slow down the convergence as that is necessary to ensure numerical stability. This is explained in appendix A.

Ignoring the higher order terms, equation 46 becomes:

$$\sum_j^n M_{ij} \delta_j^{g-1} + \frac{\tau}{u_i} \delta_i^{g-1} - \frac{\tau}{u_{i+1}} \delta_{i+1}^{g-1} = c_i - \sum_j^n M_{ij} f_j^{g-1} - \tau [\ln(f_i^{g-1}) - \ln(f_{i+1}^{g-1})] \quad (52)$$

for  $1 \leq i \leq n-1$ , and equation 36 becomes

$$\sum_j^n f_j^g = 1 \quad (53)$$

$$\sum_j^n \delta_j^{g-1} = 1 - \sum_i^n f_i^{g-1} \quad (54)$$

$$= 0 \quad (55)$$

for  $i = n$ , assuming  $\mathbf{f}^{g-1}$  is normalized.

A collorary of the definition 51 above is that at each iteration  $g$ ,  $\mathbf{f}^g$  will always deviation from the correct solution  $\mathbf{v}$  as follows:

Recalling that the LHS of equations 35 and 36 can be written as function  $F_i$ , and then vectorized into  $\mathbf{F}(\mathbf{v}) = \mathbf{0}$  as seen in equation 38,

$$\mathbf{F}(\mathbf{f}^g) = \mathbf{F}(\mathbf{f}^{g-1} + \alpha \mathbf{\delta}^{g-1}) = \mathbf{O}^{g-1} \quad (56)$$

where the deviation from the origin in codomain space  $\mathbf{O}^{g-1}$  is

$$(\mathbf{O}^{g-1})_i = \begin{cases} \sum_{j=1}^n M_{ij} (f_j^{g-1} + \alpha \delta_j^{g-1}) + \tau [\ln(f_i^{g-1} + \alpha \delta_i^{g-1}) - \ln(f_{i+1}^{g-1} + \alpha \delta_{i+1}^{g-1})] - c_i & \text{for } 1 \leq i \leq n-1 \\ \sum_{j=1}^n M_{ij} (f_j^{g-1} + \alpha \delta_j^{g-1}) - c_i & \text{for } i = n \end{cases} \quad (57)$$

Thus, for the answer to converge, we will require:

1.  $f_i^g$  to stay in the physically allowed bounds at the end, i.e.  $0 \leq f_i^g \leq 1$  as  $g \rightarrow \infty$
2.  $\mathbf{O}^g$  to be bounded,
3.  $\mathbf{O}^g$  to decrease with every iteration, converging to zero,

Condition 3 is more difficult to define, since there are  $n$  different dimensions to consider, each dimension can converge towards zero at different rates and manner (oscillatorily or monotonically); and if we use its  $p$ -norm, we can choose

$$\lim_{g \rightarrow \infty} \|\mathbf{O}\|_1 = \sum_i (f_i) \rightarrow 0, \quad \lim_{g \rightarrow \infty} \|\mathbf{O}\|_2 = \sqrt{\sum_i (f_i)^2} \rightarrow 0,$$

or any other positive values of  $p$   $\lim_{g \rightarrow \infty} \|\mathbf{O}\|_p \rightarrow 0$ , as the measurement of “convergence”.

In practice condition 3 will not be implemented as  $\mathbf{O}^g < \mathbf{O}^{g-1}$  is not guaranteed when  $\mathbf{f}$  approaches the solution and becomes

Condition 2 is easier to impose: as long as  $f_i^g > 0$ ,  $\ln(f_i^g)$ , and therefore,  $\sum_{l=2}^{\infty} \frac{1}{l} \left(\frac{\delta_i^{g-1}}{f_i^{g-1}}\right)^l$ , will not diverge to infinity, staying bounded. This can be achieved by ensuring that

$$f_i^g = f_i^{g-1} + \alpha \delta_i^{g-1} > 0 \quad \forall 1 \leq i \leq n \quad (58)$$

Given that  $\mathbf{f}_i^{DEF}$  and  $\mathbf{f}_i^g$  are both positive, condition 1 is obeyed naturally when the inequality 58 is obeyed, due to the normalization constraint imposed by equation 31.

The details of how to impose inequality 58 will be stated in appendix A.

When an appropriate choice of  $\alpha$  has been made, the procedure of finding  $\delta$  can be stated as follows: Starting from equation 52 and 54,

$$\sum_j^n M_{ij} \delta_j^{g-1} + \frac{\tau}{u_i} \delta_i^{g-1} - \frac{\tau}{u_{i+1}} \delta_{i+1}^{g-1} = c_i - \sum_j^n M_{ij} f_j^{g-1} - \tau [\ln(f_i^{g-1}) - \ln(f_{i+1}^{g-1})] \quad (59)$$

for  $1 \leq i \leq n-1$

$$\sum_j^n \delta_j^{g-1} = c_i - \sum_j^n M_{ij} f_j^{g-1} \quad (60)$$

for  $i = n$ .

To express equations 59 and 60 in vector form,

$$[\underline{\underline{\mathbf{M}}} + \tau \underline{\underline{\mathbf{A}}}^{g-1}] \boldsymbol{\delta}^{g-1} = \mathbf{c} - \underline{\underline{\mathbf{M}}} \mathbf{f}^{g-1} - \tau \underline{\underline{\mathbf{L}}}(\mathbf{f}^{g-1}) \quad (61)$$

where each of the symbol is defined as follows:

$\underline{\underline{\mathbf{M}}}$  and  $\mathbf{c}$  are as defined in equation 33 and 34,  $\tau$  (the Tikhonov parameter) is a scalar,

$$(\underline{\underline{\mathbf{L}}}(\mathbf{f}^{g-1}))_i = \begin{cases} \ln(f_i^{g-1}) - \ln(f_{i+1}^{g-1}) & \text{for } 1 \leq i \leq n-1 \\ 0 & \text{for } i = n \end{cases} \quad (62)$$

$$(\underline{\underline{\mathbf{A}}}^{g-1})_{ij} = \begin{cases} \frac{\delta_{ij}^*}{f_j^{g-1}} - \frac{\delta_{(i+1)j}^*}{f_j^{g-1}} & \text{for } 1 \leq i \leq n-1 \\ 0 & \text{for } i = n \end{cases} \quad (63)$$

where  $\delta^*$  is the Kronecka delta; the  $*$  is used to prevent clashing symbol with the vector  $\delta^{g-1}$  previously used.

The first  $n - 1$  rows of  $\underline{\underline{\mathbf{A}}}^{g-1}$  forms a double diagonal matrix (populated only at the the main diagonal and the diagonal line immediately above it).

$$\tau \underline{\underline{\mathbf{A}}}^{g-1} = \begin{pmatrix} \frac{\tau}{f_1^{g-1}} & \frac{-\tau}{f_2^{g-1}} & & & & \\ & \ddots & \ddots & & & \\ & & \frac{\tau}{f_i^{g-1}} & \frac{-\tau}{f_{i+1}^{g-1}} & & \\ & & & \ddots & \ddots & \\ & & & & \frac{\tau}{f_{n-1}^{g-1}} & \frac{-\tau}{f_n^{g-1}} \\ 0 & \dots & 0 & \dots & 0 & 0 \end{pmatrix} \quad (64)$$

Note that  $\underline{\underline{\mathbf{L}}}(\mathbf{f})$  is actually a non-linear operation on  $\mathbf{f}$ , but still outputs a vector. It can be similarly expressed as

$$\tau \underline{\underline{\mathbf{L}}} = \begin{pmatrix} \tau \ln & -\tau \ln & & & & \\ & \ddots & \ddots & & & \\ & & \tau \ln & -\tau \ln & & \\ & & & \ddots & \ddots & \\ & & & & \tau \ln & -\tau \ln \\ 0 & \dots & 0 & \dots & 0 & 0 \end{pmatrix} \quad (65)$$

The operator on the LHS of equation 61 is the Jacobian matrix,

$$\underline{\underline{\mathbf{J}}}^{g-1} = \underline{\underline{\mathbf{M}}} + \tau \underline{\underline{\mathbf{A}}}^{g-1} \quad (66)$$

#### 4.4.1 Summary

From 61,

$$\underline{\underline{\mathbf{J}}}^{g-1} \delta^{g-1} = \mathbf{c} - \underline{\underline{\mathbf{M}}} \mathbf{f}^{g-1} - \tau \underline{\underline{\mathbf{L}}}(\mathbf{f}^{g-1}) \quad (67)$$

At each step, when  $\mathbf{f}^{g-1}$  is known, the RHS is evaluated, and then the inverse of the Jacobian is left-multiplied onto both sides to obtain  $\delta^{g-1}$

$$\delta^{g-1} = \left( \underline{\underline{\mathbf{J}}}^{g-1} \right)^{-1} [\mathbf{c} - \underline{\underline{\mathbf{M}}} \mathbf{f}^{g-1} - \tau \underline{\underline{\mathbf{L}}}(\mathbf{f}^{g-1})] \quad (68)$$

And then, given the underrelaxation factor  $\alpha$ , the next iteration's  $\mathbf{f}$  can be evaluated using equation 51

$$\mathbf{f}^g = \mathbf{f}^{g-1} + \alpha \delta^{g-1} \quad (69)$$

For the first iteration  $g = 1$ , a very obvious choice of  $\mathbf{f}^{g-1} = \mathbf{f}^0$  is the *a priori*, i.e.  $\mathbf{f}^0 = \mathbf{f}^{DEF}$ . When  $\mathbf{f}^g$  settles onto an equilibrium, a check  $\mathbf{F}(\mathbf{f}^g) = \mathbf{0} \pm \text{tolerance}$  is performed. If it passes this check, then the program can terminate, returning this value as the output.

## 5 Benchmarking/Verification

## 6 Renormalization

A normalization factor  $Y$  can be added into the derivation, scaling the list of calculated reaction rates up/down, minimizing the  $\chi^2$  at the same time:

$$Y = \text{n. yield} \quad (70)$$

$$YN_k = a_k \quad (71)$$

where  $\mathbf{a}$  is the list of total reaction rates measurements at  $t=0$ , so that now  $\chi^2$  is defined as

$$\chi^2 = \sum_k^m \frac{(YN'_k - a_k)^2}{\sigma_k^2} \quad (72)$$

$$Y = \arg \min_Y (\chi^2) \quad (73)$$

This means this regularization neutron spectrum unfolding method can simultaneously output an estimate of the n. yield.

A similar procedure is carried out by MAXED at the start of each run; but this n. yield estimate is based on its comparison with the *a priori* 's reaction rates, not the final reaction rates. By imposing the constraint stated in equation 73 at each iteration, a more accurate n. yield can be obtained.

## 7 Using multiple a priori

## 8 Future direction for development

As briefly mentioned in section 3.1, the covariances between daughter nuclides resulting from the irradiation of the same foil, or from division by the n. yield diagnostics, can be accounted for by using a more sophisticated model.

Additionally, the error from the reaction rates measurements (and potentially from the cross-section values as well) can be propagated, analytically (instead of using Monte Carlo methods), through the unfolding process to obtain a covariance matrix of the flux values.

## References

- [1] J Cvachovec and F Cvachovec. Maximum likelihood estimation of a neutron spectrum and associated uncertainties. Advances in Military Technology, 1(2):67–79, 2008.
- [2] Fatma Zohra Dehimi, Abdeslam Seghour, and Saddik El Hak Abaidia. Unfolding of neutron energy spectra with fisher regularisation. IEEE Transactions on Nuclear Science, 57(2):768–774, 2010.
- [3] Shigetaka Maeda and Tetsuo Iguchi. A new unfolding code combining maximum entropy and maximum likelihood for neutron spectrum measurement. Journal of Nuclear Science and Technology, 50(4):381–386, 2013.
- [4] M. Matzke. Unfolding procedures. Radiation Protection Dosimetry, 107(1-3):155–174, 2003.

- [5] Jan Mlynar, John M. Adams, Luciano Bertalot, and Sean Conroy. First results of minimum fisher regularisation as unfolding method for jet ne213 liquid scintillator neutron spectrometry. Fusion Engineering and Design, 74(1):781 – 786, 2005. Proceedings of the 23rd Symposium of Fusion Technology.
- [6] M. Reginatto and P. Goldhagen. Maxed, a computer code for maximum entropy deconvolution of multisphere neutron spectrometer data. Health Phys., 77(5):579–83, 1999.
- [7] Marcel Reginatto, Paul Goldhagen, and Sonja Neumann. Spectrum unfolding, sensitivity analysis and propagation of uncertainties with the maximum entropy deconvolution code maxed. Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, 476(1-2):242–246, 2002.
- [8] Wolfgang von der Linden. Maximum-entropy data analysis. Applied Physics A, 60(2):155–165, 1995.

# Appendices

## A Choosing an underrelaxation factor

Consider equation 46, after the linearization approximation and removing  $O$  (higher order terms), the algorithm assumes that this is a simple linear algebra problem, i.e. the gradients in the landscapes stays constant everywhere. Therefore it takes a step as big as it think it needs to according to the gradient evaluated at  $\mathbf{u}$  when there is no underrelaxation, i.e.  $\alpha = 1$ .

To fulfill Inequality 58, i.e. no flux values overshoots into the negative flux region, we can first define the vector intended fractional increase as  $\mathbf{w}$

$$w_i = \frac{\delta_i^{g-1}}{f_i^{g-1}} \quad (74)$$

Obviously, if the intended fractional increase for the  $i^{th}$  element is less than -100% or worse, then the resulting  $f_i^g$  will definitely overshoot into the negative flux region if  $\alpha$  is kept at 1.

To moderate this, we can choose  $0 < \alpha < 1$  such that the worst fractional decrease is still less than -1, so  $O(u_i, \Delta_i)$  doesn't diverge too much.

$$x = |\min(\mathbf{w})| \quad (75)$$

$$\alpha = \min\left(\frac{1}{2x}, 1\right) \quad (76)$$

The factor of 2 in the denominator of equation 76 is chosen to minimize the amount of deviation due to the term  $O(u_i, \Delta_i)$ . A conservative algorithm will use a larger denominator (e.g. 100), sliding down the gradient smoothly and slowly, evaluating the gradient at every point along the way; while a greedy algorithm will use a smaller denominator (e.g. 1.255<sup>4</sup>) and therefore larger steps, occasionally overshoot due to underestimation of the steepness of the loss-value landscape (due to the linearization assumption), taking an oscillatory path to the correct answer.

---

<sup>4</sup>To ensure that the higher order terms in the Taylor expansion of  $\ln()$  is still smaller than the first order term,  $O(\frac{\delta}{x}) = |-\ln(\frac{\delta}{x} + 1) + 1| < |\frac{\delta}{x}|$ . Solving this inequality graphically gives  $\frac{1}{|\frac{\delta}{x}|} > \frac{1}{0.7968} \approx 1.255$ .

Note that since  $\mathbf{w}$ , which should be written as  $\mathbf{w}^{g-1}$  if one were to be scrupulous, is different across different iterations. So the maximum value that  $\alpha$  (again, should be written as  $\alpha^{g-1}$  if one were to be scrupulous) can take is also dependent on the iteration.

## B Uniqueness of global minimum

If one chooses the wrong loss function, it may lead to the existence of multiple local minima. Therefore it is important to show, under what condition will the loss function landscape have multiple local minima.

The following procedure can be applied onto any loss function with a loss function gradient composed of one linear part and one non-linear part. For our loss function  $Loss(\phi_{sol}) = (\chi^2)_{\text{reaction rates}}(\phi_{sol}) + \tau \cdot \text{RelativeEntropy}_{\phi_0}(\phi_{sol})$ ,  $\chi^2$  is the linear part, relative-entropy is the non-linear part.

The approach used in the following derivation is to think of the solution space  $\mathbf{f}$  as the domain, the gradient on the loss function  $\mathbf{G} = \nabla(loss) + \text{const.}$  as the codomain. Note that  $\nabla loss$  is a function of  $\mathbf{f}$

First, it will be beneficial to review some of the basic terminology used. Additionally, for convenient communication, I will create a term called “hyperedge”.

word	meaning	e.g. in 1D	3D	in the current context
domain	input space	x	$\mathbf{u} = (u_1, u_2, u_3)^T$	$\mathbf{f} = (f_1, \dots, f_n)^T$
codomain	output space	$y=f(x)$	$\mathbf{v} = \mathbf{F}(\mathbf{u})$	$\mathbf{G} = \nabla(loss) + \text{const.}$
Jacobian $\underline{\mathbf{J}}$	how much is the domain “stretched” to form the codomain	$\frac{dy}{dx}$	$\underline{\mathbf{J}} = \frac{\partial \mathbf{v}}{\partial \mathbf{u}}; J_{ij} = \frac{\partial v_i}{\partial u_j}$	$J_{ij} = \frac{\partial G_i}{\partial f_j}$
hyperplane	an n-1 dimensional object	a point	a plane	
hyperedge	an n-2 dimensional object	N/A	a line	
simplex	object formed by n+1 vertices	a line	a tetrahedron	

Additionally, the definition of some less-commonly seen matrix decomposition methods is provided below:

SVD : singular-value-decomposition, equivalent of eigenvalue decomposition for a singular matrix.

---

Left-singular-vector: equivalent of an eigenvector for a singular matrix.

In the following derivation, it is helpful to think of the domain and codomain as separate vector spaces. Trying to plot  $f_i$  against  $G_i$  will only confuse the reader. Instead, we should think about the transformation from  $\mathbf{f}$  into  $\mathbf{G}$ , where  $\mathbf{f}$  and  $\mathbf{G}$  are n-dimensional vector spaces. For n=2 & 3, the transformation is trivial when it’s still linear; and slightly more difficult to imagine when it’s non-linear. If you are adventurous, you can try to imagine n=4.

In section 4.1, the procedure of using Lagrangian multipliers led to equation 30 and 31. It is not difficult to see that equation 31 is a constraint that forces the solution to exist within a hyperplane in the domain space. However, perhaps it is more difficult to see that 30 is an equation that requires our solution to exist on a line  $\mathbf{c} - \mu \hat{\mathbf{n}}$  in the codomain space. This is because, if we express the LHS of 30 as  $\mathbf{G}$  (which is a function of  $\mathbf{f}$ ), and condense the RHS into a constant vector  $\mathbf{c}'$ , then we have

$$G_i(\mathbf{f}) + \mu = c'_i \quad (77)$$

So that in codomain space, we have

$$\mathbf{G} - \mathbf{c}' = -\mu \hat{\mathbf{n}}, \quad (78)$$

where  $\mu$  is an arbitrary real number, and

$$\hat{\mathbf{n}} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad (79)$$

The hyperplane in the domain space defined by equation 31 is mapped (in a non-linear manner, as dictated by the loss function's gradient) onto the codomain space,  $\mathbf{f} \rightarrow \mathbf{G}$ . The intersection points of it with the line  $\mathbf{c}' - \mu \hat{\mathbf{n}}$  are the solutions.

## B.1 Linear transformation

Now, if the mapping from  $\mathbf{f} \rightarrow \mathbf{G}$  is (almost) linear (i.e. the terms in the loss function with higher than  $2^{nd}$  order dependence on  $\mathbf{f}$  on are vanishingly small), then we can be sure that there will be:

1. no solution if line  $(\mathbf{c}' - \mu \hat{\mathbf{n}})$  does not intersect the hyperplane in the codomain space
2. 1 unique solution if the line penetrates the hyperplane at 1 point.
3. multiple (infinitely many) solutions if the line lies inside the hyperplane.

If the the hyperplane in the codomain space is parallel to the line  $(\mathbf{c}' - \mu \hat{\mathbf{n}})$ , then either situation 3 or situation 1 may occur. But the occurrence of situation 1 does not guarantee that the hyperplane is parallel to the line  $(\mathbf{c}' - \mu \hat{\mathbf{n}})$ . This is because the allowed range of solution is bounded in the domain space by  $f_i > 0$ . (In 3D, this hyperplane can be visualized as the triangle “covering” the origin, with the x-,y-, and z-intercept at 1,1, and 1.) Therefore, when a linear mapping is applied to transform this hyperplane from the domain space into the codomain space, an in appropriate choice of  $\mathbf{c}$  will lead to no solution. This example (where there is no solution even when the line  $(\mathbf{c}' - \mu \hat{\mathbf{n}})$  is not parallel to the hyperplane) corresponds to the cases where the fluxes in some bins must encroach on the negative fluxes regions in order to achieve optimal loss functions.

## B.2 when the linear transformation is singular

Another problem is that  $\underline{\underline{\mathbf{A}}}$  will be a singular matrix. Since the matrix  $\underline{\underline{\mathbf{A}}}$  is constructed from the sum of  $m$  outer products of length  $n$  vectors, the number of non-zero eigenvalues of the matrix  $\underline{\underline{\mathbf{A}}}$  is  $\leq m$ . Therefore there are at least  $(n-m)$  eigenvalues which equal to zero<sup>5</sup>. Meaning the linear transformation will map the  $\sum f_i = 1$  hyperplane in the domain into a  $m$ -dimensional manifold in the codomain. In the rare case that the  $\hat{\mathbf{n}}$  vector in the codomain is contained by this  $m$ - (or fewer-)

---

<sup>5</sup>Be careful when dealing with singular matrices such as  $\underline{\underline{\mathbf{A}}}$ , as computer programs in general will compute (incorrectly) to give non-zero eigenvalues, and successfully producing an “inverse” matrix without raising a warning.



dimensional manifold, it is possible to have infinitely many solutions, with at least (n-m) degrees of freedom in the solution space.

(Notice that the measured reaction rates only control  $\mathbf{c}$ . Therefore if the m-manifold is parallel to the  $\hat{\mathbf{n}}$  vector (i.e. has bases vectors that can span  $\hat{\mathbf{n}}$ ) in the codomain (i.e. condition 3 is met), then that means we have a response matrix where only a very specific ratio of measured reaction rates will give a solution at all. Outside of this range of physically possible measured reaction rates ratios, we will require  $\sum f_i \neq 1$  in order to fit to it optimally. The fact that  $\hat{\mathbf{n}}$  is parallel to the manifold means we have a response matrix that causes no change at all in total  $\chi^2$  even after changing the fluxes in all manner imaginable. I.e. we must have a flat cross-section in that case.)

A sensible choice of the response matrix  $\underline{\mathbf{R}}$  (i.e. one without constant cross sections) will lead to a m-manifold that does not span the vector  $\hat{\mathbf{n}}$ . However, if one is feeling paranoid, they can check that their m-manifold object in the codomain space does not span the  $\hat{\mathbf{n}}$  vector by  $\underline{\mathbf{A}}$  matrix does not by checking that all m right-singular-vectors  $\theta_k$

$$\hat{\mathbf{n}} - \sum_k^m (\hat{\mathbf{n}} \cdot \theta_k) \theta_k \neq \mathbf{0} \quad (80)$$

### B.3 Non-linear transformation

In this section we will look more closely at loss functions which have a linear part and a non-linear part.

We have designed our loss function to increase logarithmically, so harsher penalty will be given by the loss function as  $f_i \rightarrow 0^+$ , as explained in section 2. Thus, the problems discussed in section B.1 and B.2 are usually avoided when using this loss function,

$$LossFunction(\mathbf{f}, \mathbf{f}^{DEF}) = \underline{\mathbf{A}}\mathbf{f} + \tau \sum_i^n (f_i \ln(\frac{f_i}{f_i^{DEF}})) \quad (81)$$

The first term in the loss function is linear, while the second term is non-linear.

For the moment, let's assume that  $\underline{\mathbf{A}}$  is non-singular, and  $\tau$  is very small, i.e. for values of  $f_i$  far away from zero, the second term is effectively zero.

Therefore, the majority of the hyperplane in codomain space will look like  $\underline{\mathbf{A}}\mathbf{f} \{ \forall \sum_i f_i = 1 \}$ . In 3D, this is the same as the triangle formed by  $\mathbf{A}_1$ ,  $\mathbf{A}_2$ ,  $\mathbf{A}_3$ , which are the three column vectors that makes up  $\underline{\mathbf{A}}$ . Generalized to higher dimension, it is the (n-1) simplex formed by the n column vectors of  $\underline{\mathbf{A}}$ .

At this point, it is convenient to define the following vector  $\mathbf{H}$ , such that

$$\mathbf{H} = \underline{\mathbf{A}}\mathbf{f} \quad (82)$$

is the partial transformation of  $\mathbf{f}$  into codomain  $\mathbf{G}$ .

$$G_i = H_i + \tau \ln(f_i) \quad (83)$$

Note that the Jacobian of these transformations are as follows:

$$(\underline{\mathbf{J}}_{\mathbf{f} \rightarrow \mathbf{H}})_{ij} = A_{ij} \quad (84)$$

$$(\underline{\mathbf{J}}_{\mathbf{f} \rightarrow \mathbf{G}})_{ij} = A_{ij} + \tau \frac{\delta_{ij}^*}{f_i} \quad (85)$$

where  $\delta^*$  is, again, the Kronecka delta.

After each of the above transformation, the direction of the hyperplane's normal vector is given as (without proof)  $\hat{\mathbf{n}} \cdot \underline{\underline{\mathbf{J}}}^{-1}$ , where  $\hat{\mathbf{n}}$  is the unnormalized vector defined in 79. But note that this equation is only meaningful if A is non-singular, which is not true for the case when the is underdetermined.

Instead all we can do is check that the

## B.4 Other failed approaches

In science it is equally important to record what failed and why. Here are the list of other starting points that led to nowhere.

- ‘Using Stokes theorem as the starting point: leads to no useful results/intuition
- ‘Using Divergence theorem as the starting point: only gives information about the curvature of an already known point, not useful for finding points where the gradient = 0
- ‘Draw a line  $(\mathbf{u} + \lambda \mathbf{v})$ , anchored at a known global minimum  $\mathbf{u}$ , pointed at an arbitrary direction  $\mathbf{v}$ . If the loss function has a point  $\frac{\partial(loss)}{\partial \lambda} = 0$  where  $\lambda \neq 0$ , then we know that a local minimum is present’: Not feasible without knowing which directions  $\mathbf{v}$  to draw the line in in the first place.