

Protein abundances vs. equilibrium activities

Jeffrey M. Dick

June 16, 2012

1 Introduction

The `diagram()` function serves multiple purposes that might be confusing to the new user. From its name, we know that it produces diagrams of some sort. These are equilibrium chemical activity diagrams – that is the primary purpose of the function. However, inspecting the arguments to the function reveals that the input values are the affinities of formation reactions of species in the system. How do we go from chemical affinities to chemical activities? This problem defines the purpose of two auxiliary functions (`equil.react()` and `equil.boltz()`) whose algorithms are described below.

Some explanation of terminology is in order. By chemical activity we mean the quantity a_i that appears in the expression

$$\mu_i = \mu_i^\circ + RT \ln a_i, \quad (1)$$

where μ_i and μ_i° stand for the chemical potential and the standard chemical potential of the i th species, and R and T represent the gas constant and the temperature in Kelvin. Chemical activity is related to molality (m_i) by

$$a_i = \gamma_i m_i, \quad (2)$$

where γ_i stands for the activity coefficient of the i th species. For this discussion, we take $\gamma_i = 1$ for all species, so chemical activity is assumed to be numerically equivalent to molality. Since molality is a measure of concentration, calculating the equilibrium chemical activities can be a theoretical tool to help understand the relative abundances of species, including proteins.

After going over the methods used in CHNOSZ for equilibrium activity calculations, some comparisons with experimental protein abundance data are made.

2 Calculations at a single point

Here we discuss two procedures for calculating equilibrium activities of species. The first is a reaction-matrix approach and the second takes advantage of the Boltzmann distribution. We show (by example) that the two approaches are equivalent when the formation reactions of residue equivalents of proteins are used. The example system here has also been described in a paper [3].

2.1 Reaction-matrix approach

The next two sections give examples of calculating the equilibrium activities of two proteins using a matrix of equations representing reactions to form the proteins. Although the examples below include only two proteins, each additional protein introduces one more equation and unknown, so this procedure can be carried out for any number of proteins given the necessary computational requirements.

2.1.1 Whole proteins

Let us calculate the equilibrium activities of two proteins in metastable equilibrium. To do this we start by writing the formation reactions of each protein as



and



The basis species in the reactions are collectively symbolized by *stuff*; the subscripts simply refer to the reaction number in this document. In these examples, *stuff* consists of CO₂, H₂O, NH₃, O₂, H₂S and H⁺ in different molar proportions. To see what *stuff* is, try out these commands in CHNOSZ:

```
> library(CHNOSZ)
> basis("CHNOS+")
```

	C	H	N	O	S	Z	ispecies	logact	state
CO2	1	0	0	2	0	0	69	-3	aq
H2O	0	2	0	1	0	0	1	0	liq
NH3	0	3	1	0	0	0	68	-4	aq
H2S	0	2	0	0	1	0	70	-7	aq
O2	0	0	0	2	0	0	3097	-80	gas
H+	0	1	0	0	0	1	3	-7	aq

```
> species("CSG",c("METVO","METJA"))
```

	CO2	H2O	NH3	H2S		O2	H+	ispecies	logact	state	name
1	2575	1070	645	11	-2668.0	0		3371	-3	aq	CSG_METVO
2	2555	1042	640	14	-2643.5	0		3372	-3	aq	CSG_METJA

Although the basis species are defined, the temperature is not yet specified, so it is not immediately possible to calculate the ionization states of the proteins. That is why the coefficient on H⁺ is zero in the output above. To see what the computed protein charges are at 25 °C and 1 bar and at pH 7 (which is the opposite of the logarithm of activity of H⁺ in the basis species), try this:

```
> protein.info(species())$name)
```

```
subcrt: 2 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)
```

	protein	length		formula	G
1	CSG_METVO	553	C2575H4040.93490596228N6450884S11-56.0650940377185	-24406217	
2	CSG_METJA	530	C2555H3976.12975396577N6400865S14-55.8702460342319	-23718389	
	Z	G.Z	ZC		
1	-56.06509	-24502046	-0.1444660		
2	-55.87025	-23895850	-0.1385519		

Note that `affinity()` is called twice by `protein.info()`; this so that both charges and standard Gibbs energies of ionization of the proteins can be calculated. The Z values in the table are the charges of the proteins computed using the ionization constants of sidechain and terminal groups, and the G.Z values are the calculated Gibbs energies of formation of the ionized proteins [4]. The ZC values are the average oxidation states of carbon of the proteins. Let us now calculate the chemical affinities of formation of the ionized proteins:

```
> a <- affinity()
```

```
energy.args: temperature is 25 C
energy.args: pressure is Psat
subcrt: 8 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)
```

```
> a$values
```

```
$`3371`  
[1] -240.2934
```

```
$`3372`  
[1] -62.41683
```

Since `affinity()` returns a list with a lot of information (such as the basis species and species definitions) the last command was written to only print the values part of that list. The values are actually dimensionless, i.e. $A/2.303RT$.

The affinities of the formation reactions above were calculated for a *reference value of activity of the proteins, which is not the equilibrium value*. Those non-equilibrium activities were 10^{-3} . How do we calculate the equilibrium values? Let us write specific statements of the expression for chemical affinity (2.303 is used here to stand for $\ln 10$),

$$A = 2.303RT \log(K/Q), \quad (5)$$

for Reactions 3 and 4 as

$$\begin{aligned} A_3/2.303RT &= \log K_3 - \log Q_3 \\ &= \log K_3 + \log a_{stuff,3} - \log a_{CSG_METVO} \\ &= A_3^*/2.303RT - \log a_{CSG_METVO} \end{aligned} \quad (6)$$

and

$$\begin{aligned} A_4/2.303RT &= \log K_4 - \log Q_4 \\ &= \log K_4 + \log a_{stuff,4} - \log a_{CSG_METJA} \\ &= A_4^*/2.303RT - \log a_{CSG_METJA}. \end{aligned} \quad (7)$$

The A^* denote the affinities of the formation reactions when the activities of the proteins are zero. From the output above it follows that $A_3^*/2.303RT = 104.6774$ and $A_4^*/2.303RT = 314.1877$.

Next we must specify how reactions are balanced in this system: what is conserved during transformations between species (let us call it the immobile component)? For proteins, one possibility is to use the repeating protein backbone group. Let us use n_i to designate the number of residues in the i th protein, which is equal to the number of backbone groups, which is equal to the length of the sequence. If $\gamma_i = 1$ in Eq. (2), the relationship between the activity of the i th protein (a_i) and the activity of the residue equivalent of the i th protein ($a_{residue,i}$) is

$$a_{residue,i} = n_i \times a_i. \quad (8)$$

We can use this to write a statement of mass balance:

$$553 \times a_{CSG_METVO} + 530 \times a_{CSG_METJA} = 1.083. \quad (9)$$

At equilibrium, the affinities of the formation reactions, per conserved quantity (in this case protein backbone groups) are equal. Therefore $A = A_3/553 = A_4/530$ is a condition for equilibrium. Combining this with Eqs. (6) and (7) gives

$$A/2.303RT = (104.6774 - \log a_{CSG_METVO}) / 553 \quad (10)$$

and

$$A/2.303RT = (314.1877 - \log a_{CSG_METJA}) / 530. \quad (11)$$

Now we have three equations (9–11) with three unknowns. The solution can be displayed in CHNOSZ as follows. The argument `residue=FALSE` overrides the default setting for `diagram` when proteins are the species of interest and instructs it to use the function named `equil.react()`, which implements the equation-solving strategy described in the next section. Here we retrieve the equilibrium activities using `diagram()` without letting it actually do any plotting.

```
> d <- diagram(a,residue=FALSE,plot.it=FALSE)
```

diagram: balanced quantity is moles of protein backbone group
 diagram: balancing coefficients are 553 530
 diagram: log total activity of PBB (from species) is 0.0346284566253204

```

> d$logact

$`3371`
[1] -177.8441

$`3372`
[1] -2.689647

```

Those are the logarithms of the equilibrium activities of the proteins. Combining these values with either Eqs. (10) or (11) gives us the same value for affinity of the formation reactions per residue (or per protein backbone group), $A/2.303RT = 0.5978817$. Equilibrium activities that differ by such great magnitudes make it appear that the proteins are very unlikely to coexist in metastable equilibrium. Later we explain the concept of using residue equivalents of the proteins to achieve a different result.

2.1.2 Implementing the reaction-matrix approach

The implementation used in CHNOSZ for finding a solution to the system of equations relies on a difference function for the activity of the immobile component. The steps to obtain this difference function are:

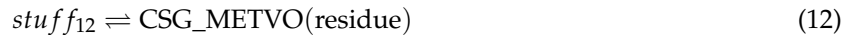
1. Set the total activity of the immobile (conserved) component as a_{ic} (e.g., the 1.083 in Eqn. 9).
2. Write a function for the logarithm of activity of each of the species of interest: $A = (A_i^* - 2.303RT \log a_i) / n_{ic,i}$, where $n_{ic,i}$ stands for the number of moles of the immobile component that react in the formation of one mole of the i th species. (e.g., for systems of proteins where the backbone group is conserved, $n_{ic,i}$ is the same as n_i in Eq. 8). Calculate values for each of the A_i^* . Metastable equilibrium is implied by the identity of A in all of the equations.
3. Write a function for the total activity of the immobile component: $a'_{ic} = \sum n_{ic,i} a_i$.
4. The difference function is now $\delta a_{ic} = a'_{ic} - a_{ic}$.

Now all we have to do is solve for the value of A where $\delta a_{ic} = 0$. This is achieved in the code by first looking for a range of values of A where at one end $\delta a_{ic} < 0$ and at the other end $\delta a_{ic} > 0$, then using the `uniroot()` function that is part of R to find the solution.

This approach is subject to failure if for all trial ranges of A the δa_{ic} are of the same sign, which gives an error message like “i tried it 1000 times but can’t make it work”. Even if values of δa_{ic} on either side of zero can be located, the algorithm does not guarantee an accurate solution and may give a warning about poor convergence if a certain (currently hard-coded) tolerance is not reached.

2.1.3 Residue equivalents

Let us consider the formation reactions of the *residue equivalents* of proteins, for example



and



The formulas of the residue equivalents are those of the proteins divided by the number of residues in each protein. With the `protein.basis()` function it is possible to see the coefficients on the basis species in these reactions:

```
> protein.basis(species())$name, residue=TRUE)
```

```

subcrt: 18 species at 298.15 K and 1 bar (wet)
      CO2      H2O      NH3      H2S      O2      H+
[1,] 4.656420 1.934901 1.166365 0.01989150 -4.824593 -0.1013835
[2,] 4.820755 1.966038 1.207547 0.02641509 -4.987736 -0.1054156

```

Let us denote by A_{12} and A_{13} the chemical affinities of Reactions 12 and 13. We can write

$$A_{12}/2.303RT = \log K_{12} + \log a_{stuff,12} - \log a_{\text{CSG_METVO}(\text{residue})} \quad (14)$$

and

$$A_{13}/2.303RT = \log K_{13} + \log a_{stuff,13} - \log a_{\text{CSG_METJA}(\text{residue})}, \quad (15)$$

For metastable equilibrium we have $A_{12}/1 = A_{13}/1$. The 1's in the denominators are there as a reminder that we are still conserving residues, and that each reaction now is written for the formation of a single residue equivalent. So, let us write A for $A_{12} = A_{13}$ and also define $A_{12}^* = A_{12} + 2.303RT \log a_{\text{CSG_METVO}(\text{residue})}$ and $A_{13}^* = A_{13} + 2.303RT \log a_{\text{CSG_METJA}(\text{residue})}$. At the same temperature, pressure and activities of basis species and proteins as shown in the previous section, we can write $A_{12}^* = A_3^*/553 = 2.303RT \times 0.1892901$ and $A_{13}^* = A_4^*/530 = 2.303RT \times 0.5928069$ to give

$$A/2.303RT = 0.1892901 - \log a_{\text{CSG_METVO}(\text{residue})} \quad (16)$$

and

$$A/2.303RT = 0.5928069 - \log a_{\text{CSG_METJA}(\text{residue})}, \quad (17)$$

which are equivalent to Equations 12 and 13 in the paper [3] but with more decimal places shown. A third equation arises from the conservation of amino acid residues:

$$a_{\text{CSG_METVO}(\text{residue})} + a_{\text{CSG_METJA}(\text{residue})} = 1.083. \quad (18)$$

The solution to these equations is $a_{\text{CSG_METVO}(\text{residue})} = 0.3065982$, $a_{\text{CSG_METJA}(\text{residue})} = 0.7764018$ and $A/2.303RT = 0.7027204$.

The corresponding logarithms of activities of the proteins are $\log(0.307/553) = -3.256$ and $\log(0.776/530) = -2.834$. These activities of the proteins are much closer to each other than those calculated using formation reactions for whole protein formulas, so this result seems more compatible with the actual coexistence of proteins in nature.

The approach just described is not used by `diagram()` when `residue=TRUE` (which is the default setting). Instead, the Boltzmann distribution, described next, is implemented for that situation.

2.2 Boltzmann distribution

An expression for Boltzmann distribution, relating equilibrium activities of species to the affinities of their formation reactions, can be written as (using the same definitions of the symbols above)

$$\frac{a_i}{\sum a_i} = \frac{e^{A_i^*/RT}}{\sum e^{A_i^*/RT}}. \quad (19)$$

Using this equation, we can very quickly (without setting up a system of equations) calculate the equilibrium activities of proteins using their residue equivalents. Above, we saw $A_{12}^*/2.303RT = 0.1892901$ and $A_{13}^*/2.303RT = 0.5928069$. Multiplying by $\ln 10 = 2.302585$ gives $A_{12}^*/RT = 0.4358565$ and $A_{13}^*/RT = 1.364988$. We then have $e^{A_{12}^*/RT} = 1.546287$ and $e^{A_{13}^*/RT} = 3.915678$. This gives us $\sum e^{A_i^*/RT} = 5.461965$, $a_{12}/\sum a_i = 0.2831009$ and $a_{13}/\sum a_i = 0.7168991$. Since $\sum a_i = 1.083$, we arrive at $a_{12} = 0.3065982$ and $a_{13} = 0.7764018$, the same result as above. This example was also described in a recent paper [5].

This computation can be carried out in CHNOSZ using the following commands, which implies `residue=TRUE` as the default setting for systems of proteins. This setting signifies to consider the formation reactions of the residue equivalents instead of the whole proteins, AND consequently to make a call to `equil.boltz()` rather than `equil.react()`.

```
> d <- diagram(a,plot.it=FALSE)

diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 553 530
diagram: using residue equivalents
diagram: log total activity of PBB (from species) is 0.0346284566253204

> as.numeric(d$logact)

[1] -3.195608 -2.860635
```

We can also specify `as.residue=TRUE` (which means to return the logarithms of activities of the residue equivalents rather than converting them to logarithms of activities of the proteins):

```
> d <- diagram(a,as.residue=TRUE,plot.it=FALSE)

diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 553 530
diagram: using residue equivalents
diagram: log total activity of PBB (from species) is 0.0346284566253204

> 10^as.numeric(d$logact)

[1] 0.3524656 0.7305344
```

Although this example includes only two proteins, this procedure is suitable for calculating the metastable equilibrium activities of any number of proteins.

3 Calculations as a function of a single variable

A comparison of the outcomes of equilibrium calculations that do and do not use the residue equivalents for proteins was given in a publication [3]. An expanded version of a diagram in that paper is below (though, without labels on the figures).

```
> organisms <- c("METSC","METJA","METFE","HALJP","METVO","METBU","ACEKI","GEOSE","BACLI","AERSA")
> proteins <- c(rep("CSG",6),rep("SLAP",4))
> basis("CHNOS+")

  C  H  N  O  S  Z ispecies logact state
CO2 1 0 0 2 0 0      69      -3    aq
H2O 0 2 0 1 0 0       1       0    liq
NH3 0 3 1 0 0 0      68      -4    aq
H2S 0 2 0 0 1 0      70      -7    aq
O2   0 0 0 2 0 0    3097     -80    gas
H+   0 1 0 0 0 1       3      -7    aq

> species(proteins,organisms)

  CO2  H2O  NH3  H2S      O2  H+ ispecies logact state      name
1  2812 1066  747  16 -2909.0  0   3373     -3    aq  CSG_METSC
2  2555 1042  640  14 -2643.5  0   3372     -3    aq  CSG_METJA
3  2815 1071  747  14 -2914.5  0   3374     -3    aq  CSG_METFE
4  3669 1367  971   0 -3608.5  0   3375     -3    aq  CSG_HALJP
5  2575 1070  645  11 -2668.0  0   3371     -3    aq  CSG_METVO
6  1362  519  355   4 -1400.5  0   3376     -3    aq  CSG_METBU
7  3584 1431  926   4 -3730.5  0   3377     -3    aq  SLAP_ACEKI
8  5676 2320 1489   3 -5904.5  0   3378     -3    aq  SLAP_GEOSE
9  3977 1594 1068   2 -4131.0  0   3379     -3    aq  SLAP_BACLI
10 2250  861  618   2 -2322.5  0   3380     -3    aq  SLAP_AERSA
```

```

> a <- affinity(O2=c(-100,-65))

energy.args: temperature is 25 C
energy.args: pressure is Psat
energy.args: variable 1 is log_f(O2) at 128 values from -100 to -65
subcrt: 16 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)

> par(mfrow=c(2,1))
> diagram(a,ylim=c(-5,-1),legend.x=NULL,residue=FALSE)

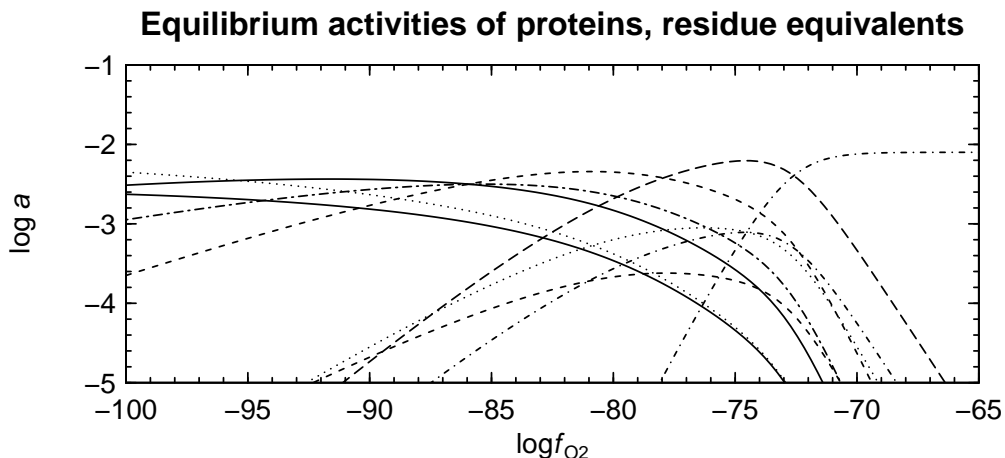
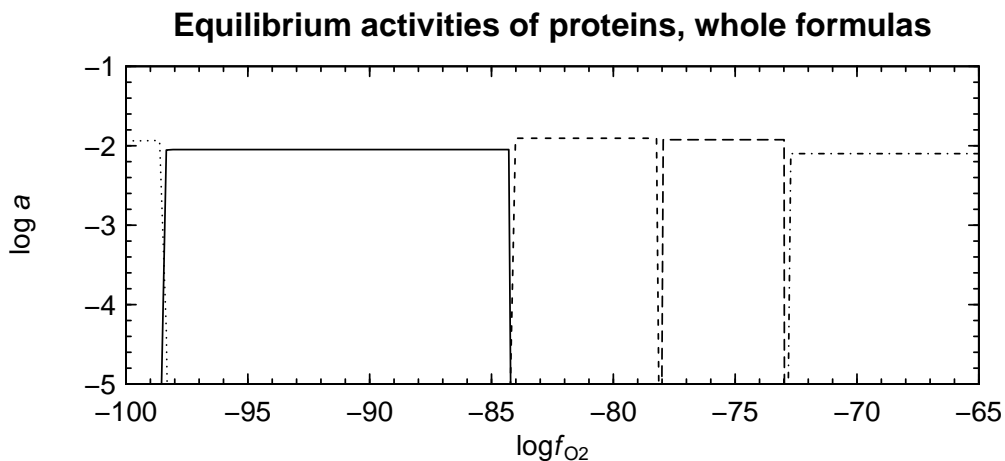
diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 571 530 571 828 553 278 736 1198 844 481
diagram: log total activity of PBB (from species) is 0.81888541459401

> title(main="Equilibrium activities of proteins, whole formulas")
> diagram(a,ylim=c(-5,-1),legend.x=NULL)

diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 571 530 571 828 553 278 736 1198 844 481
diagram: using residue equivalents
diagram: log total activity of PBB (from species) is 0.81888541459401

> title(main="Equilibrium activities of proteins, residue equivalents")

```



The reaction-matrix approach described above can also be applied to systems having conservation coefficients that differ from unity, such as many mineral and inorganic systems, where the immobile component has different molar coefficients in the formulas. For example, consider a system like that described by Seewald, 1997 [7]:

```
> basis('CHNOS+')
```

	C	H	N	O	S	Z	ispecies	logact	state
CO2	1	0	0	2	0	0	69	-3	aq
H2O	0	2	0	1	0	0	1	0	liq
NH3	0	3	1	0	0	0	68	-4	aq
H2S	0	2	0	0	1	0	70	-7	aq
O2	0	0	0	2	0	0	3097	-80	gas
H+	0	1	0	0	0	1	3	-7	aq

```
> basis('pH',5)
```

	C	H	N	O	S	Z	ispecies	logact	state
CO2	1	0	0	2	0	0	69	-3	aq
H2O	0	2	0	1	0	0	1	0	liq
NH3	0	3	1	0	0	0	68	-4	aq
H2S	0	2	0	0	1	0	70	-7	aq
O2	0	0	0	2	0	0	3097	-80	gas
H+	0	1	0	0	0	1	3	-5	aq

```
> species(c('H2S','S2-2','S3-2','S203-2','S204-2','S306-2','S506-2','S206-2','HS03-','S02','HS04-'))
```

	CO2	H2O	NH3	H2S	O2	H+	ispecies	logact	state	name
1	0	0	0	1	0.0	0	70	-3	aq	H2S
2	0	-1	0	2	0.5	-2	53	-3	aq	S2-2
3	0	-2	0	3	1.0	-2	54	-3	aq	S3-2
4	0	-1	0	2	2.0	-2	26	-3	aq	S203-2
5	0	-1	0	2	2.5	-2	1072	-3	aq	S204-2
6	0	-2	0	3	4.0	-2	1077	-3	aq	S306-2
7	0	-4	0	5	5.0	-2	1079	-3	aq	S506-2
8	0	-1	0	2	3.5	-2	1076	-3	aq	S206-2
9	0	0	0	1	1.5	-1	23	-3	aq	HS03-
10	0	-1	0	1	1.5	0	78	-3	aq	S02
11	0	0	0	1	2.0	-1	25	-3	aq	HS04-

```
> a <- affinity(O2=c(-50,-15),T=325,P=350)
```

```
energy.args: temperature is 325 C
```

```
energy.args: pressure is 350 bar
```

```
energy.args: variable 1 is log_f(O2) at 128 values from -50 to -15
```

```
subcrt: 17 species at 598.15 K and 350 bar (wet)
```

```
> par(mfrow=c(2,1))
```

```
> diagram(a,loga.balance=-2,ylim=c(-30,0),legend.x="topleft",cex.names=0.8)
```

```
diagram: balanced quantity is moles of H2S
```

```
diagram: balancing coefficients are 1 2 3 2 2 3 5 2 1 1 1
```

```
diagram: log total activity of H2S (from argument) is -2
```

```
> title(main="Aqueous sulfur speciation, whole formulas")
```

```
> diagram(a,loga.balance=-2,ylim=c(-30,0),legend.x="topleft",cex.names=0.8,residue=TRUE)
```

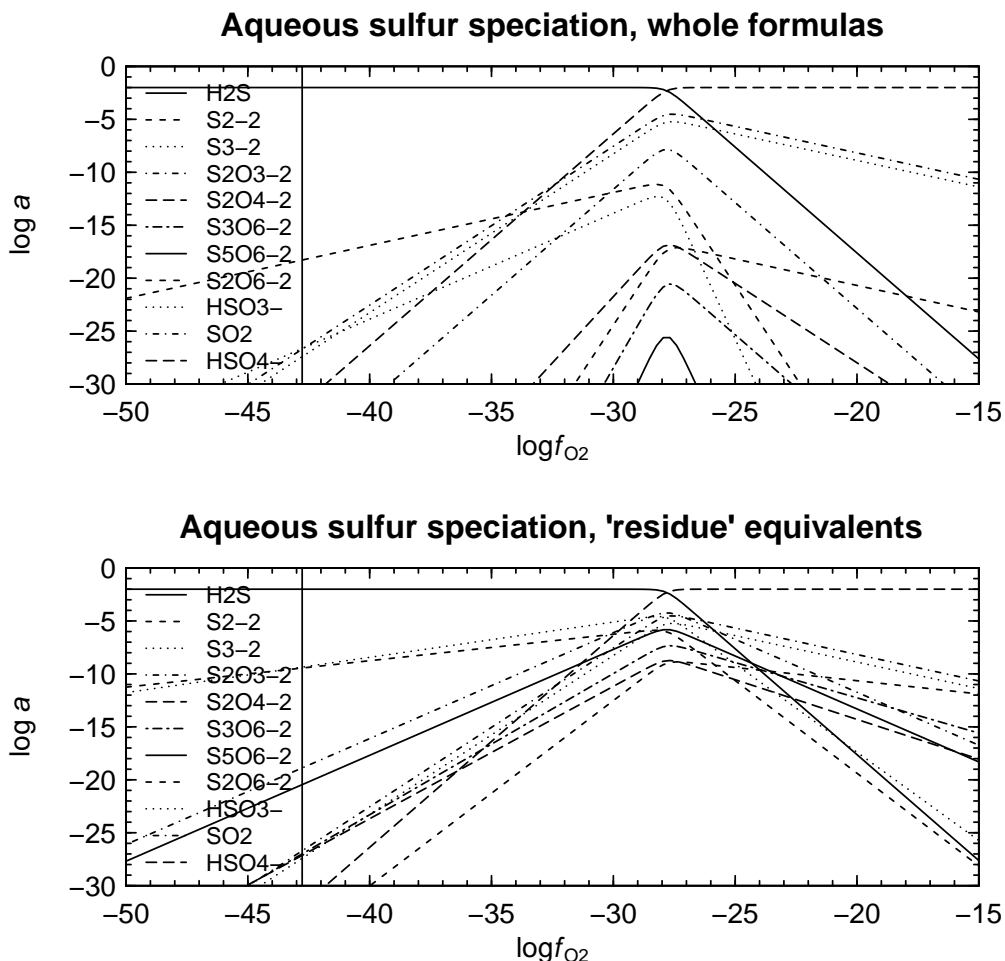


```

diagram: balanced quantity is moles of H2S
diagram: balancing coefficients are 1 2 3 2 2 3 5 2 1 1 1
diagram: using residue equivalents
diagram: log total activity of H2S (from argument) is -2

> title(main="Aqueous sulfur speciation, 'residue' equivalents")

```



The first diagram is quantitatively similar to the one shown by Seewald, 1997, but in the second (where we have set `residue=TRUE`) the range of activities is lower at any given $\log f_{\text{O}_2(g)}$. There, the function was told to rewrite the formation reactions of the aqueous sulfur species for their residue equivalents in the same way the formation reactions for the proteins were rewritten above. The number of “residues” in each species is the coefficient of the immobile component, in this case H_2S , in the formation reaction.

Maybe `residue=TRUE` doesn’t make sense for systems where the formulas of species are similar in size to those of the basis species. For molecules as large as proteins it might be a useful concept. It is now (since CHNOSZ version 0.9) the default mode for `diagram()` when working with proteins.

With the potential for calculating equilibrium activities of proteins comes the desire to compare these calculations to actual measurements!

4 Blood plasma proteins

Let’s look at some protein abundance levels in human blood plasma. First get going with the experimental data. In CHNOSZ is a table listing the upper limits of the intervals, or ranges, of protein abundances taken

from figures available in Anderson and Anderson, 2002, 2003 [1, 2]. The protein abundances in the tables are in $\log_{10}(\text{pg/ml})$; let's convert that to molality. First locate the file with the abundance data. Then read it, with "as.is" so that strings are read as characters not factors (to avoid an error in the `species()` call further down). Then identify the protein named "INS.C" and drop it from the list. The reason for doing so is that preliminary calculations show it is much more stable than any other protein in the list. It is therefore an interesting outlier in terms of relative stabilities of the proteins.

Then get the species indices of the proteins for thermodynamic calculations (with parameters based on amino acid compositions of the proteins listed in `thermo$protein ...` and suppress messages that would fill up a whole page here). Then calculate the masses of the proteins. Then convert $\log_{10}(\text{pg/ml})$ to $\log_{10}(\text{mol/L})$ (logarithm of molarity). The conversion from pg/ml to g/L involves a factor of $10^{-9} \left(\frac{10^0 \text{g}}{10^{12} \text{pg}} \times \frac{10^3 \text{ml}}{10^0 \text{L}} \right)$; then to get molarity we divide by mass ($\frac{\text{g}}{\text{mol}}$).

```
> f <- system.file("extdata/abundance/AA03.csv", package="CHNOSZ")
> pdata <- read.csv(f, as.is=TRUE)
> pdrop <- which(pdata$name == "INS.C")
> pname <- pdata$name[-pdrop]
> iip <- suppressMessages(info(paste(pname, "HUMAN", sep="_")))
> pmass <- mass(thermo$obigt$formula[iip])
> loga.expt <- logm <- log10( 10~pdata$log10.pg.ml.[-pdrop] / 10~9 / pmass )
```

As implied by the "loga", we are assuming for the comparisons offered below that molarity (derived from the published abundance data) can be taken to be equal to molality and that molality can equated with chemical activity. The latter equality (the assumption of ideal behavior) especially should be subject to more scrutiny. We'll go ahead anyway and calculate, for ideality, the equilibrium activities of the proteins. First we need to calculate the total activity of residues from the experimental data, but to do that we need even more firstly the lengths of the proteins.

```
> pl <- protein.length(paste(pname, "HUMAN", sep="_"))
> logares.tot <- sum(10~loga.expt * pl)
```

Our total activity (*not* the logarithm of it) of residues turns out to be about 200, which for our average protein length of 637 works out to about 0.3 for the average protein, if the total activity of residues could be attributed to that single average protein.

Now let's get down to the stuff CHNOSZ is made for. First define the basis species. Then define the species, being the proteins. Then calculate the affinities of the formation reactions of the proteins. Then calculate the equilibrium activities, but don't plot them by themselves. Instead, use `revisit` to compare the equilibrium activities to the experimental abundances.

```
> basis("CHNOS+")
```

	C	H	N	O	S	Z	ispecies	logact	state
C02	1	0	0	2	0	0	69	-3	aq
H2O	0	2	0	1	0	0	1	0	liq
NH3	0	3	1	0	0	0	68	-4	aq
H2S	0	2	0	0	1	0	70	-7	aq
O2	0	0	0	2	0	0	3097	-80	gas
H+	0	1	0	0	0	1	3	-7	aq

```
> s <- species(pname, "HUMAN")
> a <- affinity()
```

```
energy.args: temperature is 25 C
energy.args: pressure is Psat
subcrt: 76 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)
```

```

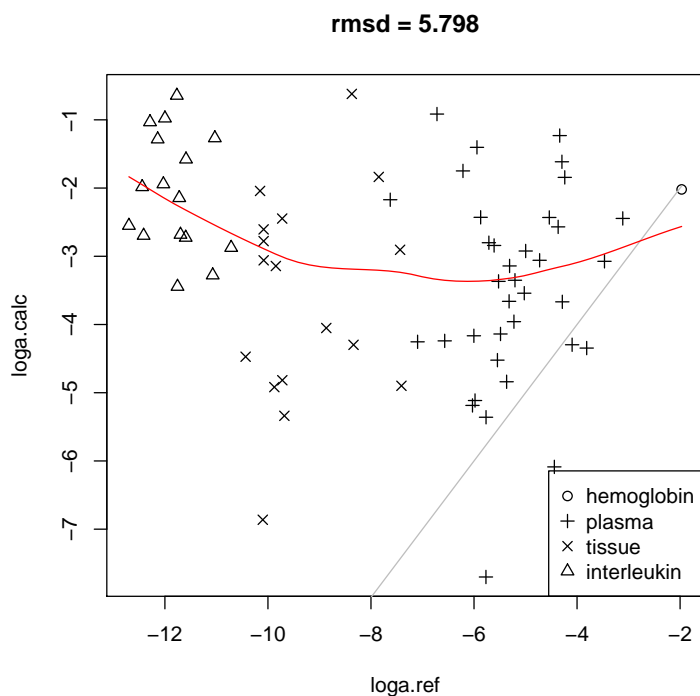
> d <- diagram(a,loga.balance=logares.tot,plot.it=FALSE)
diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 141 585 330 679 831 353 1451 137 394 1641 388 243 4536 183 4529 1213
diagram: using residue equivalents
diagram: log total activity of PBB (from argument) is 2.34870631519941

> pch <- as.numeric(as.factor(pdata$type))
> revisit(d,"rmsd",loga.ref=loga.expt,pch=pch)

revisit: calculating rmsd in 0 dimensions

> legend("bottomright",pch=unique(pch),legend=unique(pdata$type))

```



There seems to be almost no relation between the reference values and the calculated ones. But what if we increase the oxygen fugacity? $\log f_{\text{O}_{2(g)}} = -80$ might be appropriate for some subcellular conditions, or reduced hydrothermal systems. Blood is exposed to oxygen after all... let's try $\log f_{\text{O}_{2(g)}} = -60$.

```

> basis("O2",-60)

  C  H  N  O  S  Z ispecies logact state
CO2 1  0  0  2  0  0      69    -3   aq
H2O 0  2  0  1  0  0       1     0  liq
NH3 0  3  1  0  0  0      68    -4   aq
H2S 0  2  0  0  1  0      70    -7   aq
O2   0  0  0  2  0  0    3097   -60  gas
H+   0  1  0  0  0  1       3    -7   aq

> a <- affinity()

energy.args: temperature is 25 C
energy.args: pressure is Psat
subcrt: 76 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)

```

```

> d <- diagram(a,loga.balance=logares.tot,plot.it=FALSE)

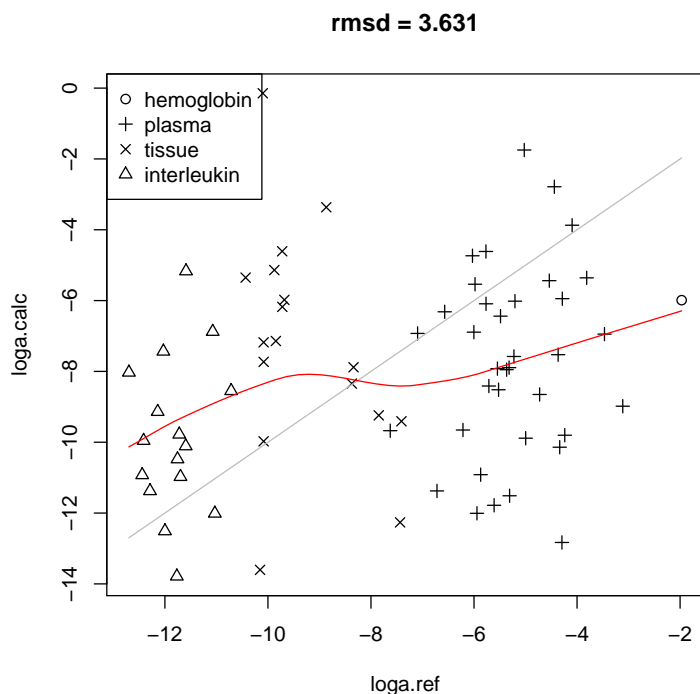
diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 141 585 330 679 831 353 1451 137 394 1641 388 243 4536 183 4529 1213
diagram: using residue equivalents
diagram: log total activity of PBB (from argument) is 2.34870631519941

> revisit(d,"rmsd",loga.ref=loga.expt,pch=pch)

revisit: calculating rmsd in 0 dimensions

> legend("topleft",pch=unique(pch),legend=unique(pdata$type))

```



Well it's still quite scattered. However, the RMSD has decreased considerably, the loess fit has a positive slope, and the dynamic ranges of the calculations and observations are more similar.

5 Comparison with expression profile in *E. coli*

Amino acid compositions of proteins in *Escherichia coli* are provided with CHNOSZ at `extdata/protein/ECO.csv.xz`. Protein abundances in the cytosol of *E. coli* reported by Ishihama et al., 2008 [6] are provided with CHNOSZ at `extdata/abundances/ISR+08.csv.xz`. We can use `read.expr()` to retrieve the abundance data for all or only selected proteins, and also add these proteins to CHNOSZ's inventory (`thermo$protein`) based on amino acid compositions from the `ECO.csv` file. First though we use `data(thermo)` to clear out the settings from the previous calculations.

```

> data(thermo)
> file <- system.file("extdata/abundance/ISR+08.csv.xz",package="CHNOSZ")
> expr <- read.expr(file,"ID","emPAI",list(description="kinase"))
> range(expr$abundance)

[1] 1.25e-01 3.38e+04

```

The result (expr) lists data for proteins where the description column of ISR+08.csv contains the term kinase. The list has elements named protein (names of proteins from the ID column of ISR+08.csv) and abundance (abundance of the proteins taken from the emPAI column of ISR+08.csv). The minimum and maximum values of the reference abundances are separated by over five orders of magnitude. Now let's use more.aa() to get the amino acid compositions of the proteins.

```
> aa <- more.aa(expr$protein, "Eco")

more.aa: KPY1 PPCK K6P1 KPY2 K6P2 were not matched

> ina <- is.na(aa$chains)
> ip <- add.protein(aa)[!ina]

add.protein: added 36 of 36 proteins
```

Note that the ID's of five of the 36 proteins that are described as "kinase" are not found in ECO.csv, so only 31 proteins are returned by the above call to more.aa().

Now let's calculate the metastable equilibrium activities of the proteins, setting the total activity of residues to unity. We then use revisit() to make a plot and compute the root mean square deviation between the experimental and calculated relative abundances. Since the equilibrium activities of the proteins were only calculated at a single point, revisit() here makes a scatter plot. The colors reflect the average oxidation state of carbon of the proteins (red – more reduced, blue – more oxidized). A crucial step here is the penultimate line, that first takes the logarithms of the observed abundances, then scales them using unitize() so that the total activity of residues is unity; this is equal to the total activity of residues in the equilibrium calculation, indicated by loga.balance=0 in diagram().

```
> basis("CHNOSZ")

  C H N O S Z ispecies logact state
CO2 1 0 0 2 0 0      69     -3   aq
H2O 0 2 0 1 0 0       1      0   liq
NH3 0 3 1 0 0 0      68     -4   aq
H2S 0 2 0 0 1 0      70     -7   aq
O2   0 0 0 2 0 0    3097    -80   gas
H+   0 1 0 0 0 1       3     -7   aq

> a <- affinity(iprotein=ip)

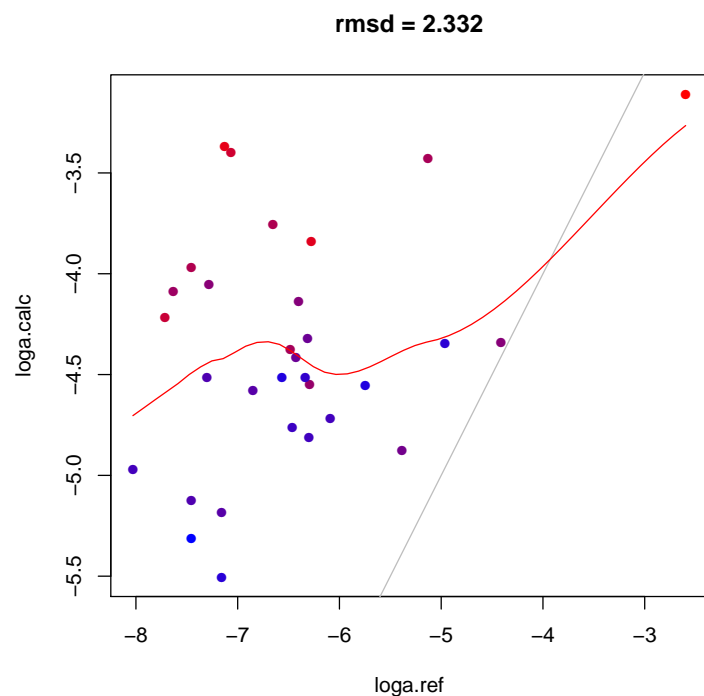
energy.args: temperature is 25 C
energy.args: pressure is Psat
subcrt: 27 species at 298.15 K and 1 bar (wet)
subcrt: 18 species at 298.15 K and 1 bar (wet)

> d <- diagram(a,loga.balance=0,plot.it=FALSE)

diagram: balanced quantity is moles of protein backbone group
diagram: balancing coefficients are 387 400 820 502 214 315 420 566 173 143 347 241 382 227 367 207 310
diagram: using residue equivalents
diagram: log total activity of PBB (from argument) is 0

> z <- ZC(protein.formula(ip))
> col <- rgb(max(z)-z, 0, z-min(z), max=diff(range(z)))
> loga.ref <- unitize(log10(expr$abundance[!ina]), length=protein.length(ip))
> revisit(d,"rmsd",loga.ref=loga.ref,pch=16,col=col)

revisit: calculating rmsd in 0 dimensions
```



How can the correlation be improved? We can find where the RMSD minimizes as a function of a single variable. Or let's go for two variables ... note that we have to specify `mam=FALSE` in the call to `diagram()` in this case:

```
> a <- affinity(O2=c(-90,-60),NH3=c(-35,0),iprotein=ip)
```

```
energy.args: temperature is 25 C
```

```
energy.args: pressure is Psat
```

```
energy.args: variable 1 is log_f(O2) at 128 values from -90 to -60
```

```
energy.args: variable 2 is log_a(NH3) at 128 values from -35 to 0
```

```
subcrt: 27 species at 298.15 K and 1 bar (wet)
```

```
subcrt: 18 species at 298.15 K and 1 bar (wet)
```

```
> d <- diagram(a,loga.balance=0,plot.it=FALSE,mam=FALSE)
```

```
diagram: balanced quantity is moles of protein backbone group
```

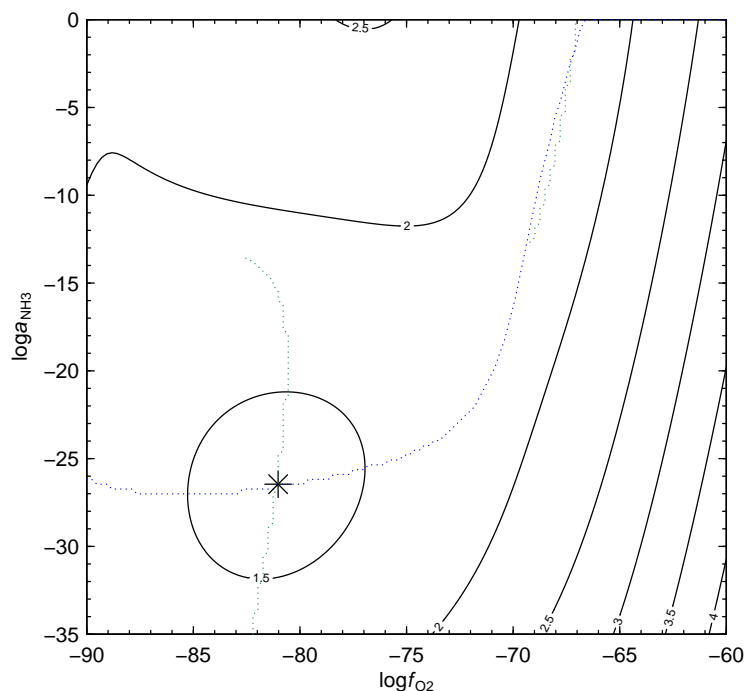
```
diagram: balancing coefficients are 387 400 820 502 214 315 420 566 173 143 347 241 382 227 367 207 310
```

```
diagram: using residue equivalents
```

```
diagram: log total activity of PBB (from argument) is 0
```

```
> r <- revisit(d,"rmsd",loga.ref=loga.ref)
```

```
revisit: calculating rmsd in 2 dimensions
```



Now set the activities of the basis species where the minimum RMSD was found, calculate the affinities and equilibrium activities, and compare the results with the reference abundances.

```
> basis(c("O2", "NH3"), c(r$x, r$y))
```

	C	H	N	O	S	Z	ispecies	logact	state
CO2	1	0	0	2	0	0	69	-3.00000	aq
H2O	0	2	0	1	0	0	1	0.00000	liq
NH3	0	3	1	0	0	0	68	-26.45669	aq
H2S	0	2	0	0	1	0	70	-7.00000	aq
O2	0	0	0	2	0	0	3097	-81.02362	gas
H+	0	1	0	0	0	1	3	-7.00000	aq

```
> a <- affinity(iprotein=ip)
```

```
energy.args: temperature is 25 C
```

```
energy.args: pressure is Psat
```

```
subcrt: 27 species at 298.15 K and 1 bar (wet)
```

```
subcrt: 18 species at 298.15 K and 1 bar (wet)
```

```
> d <- diagram(a, loga.balance=0, plot.it=FALSE)
```

```
diagram: balanced quantity is moles of protein backbone group
```

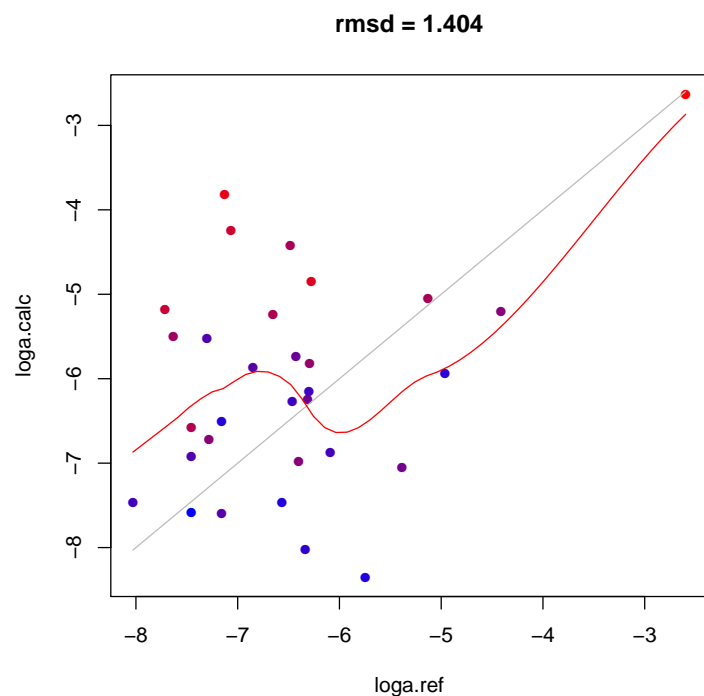
```
diagram: balancing coefficients are 387 400 820 502 214 315 420 566 173 143 347 241 382 227 367 207 310
```

```
diagram: using residue equivalents
```

```
diagram: log total activity of PBB (from argument) is 0
```

```
> revisit(d, "rmsd", loga.ref=loga.ref, pch=16, col=col)
```

```
revisit: calculating rmsd in 0 dimensions
```



Looks a little better ... still more work to do!

6 Summary

Using default settings, equilibrium activities of proteins are calculated in CHNOSZ by converting formation reactions of proteins to their per-residue equivalents, then using the Boltzmann distribution to transform the affinities of the formation reactions (in an equal-activity reference state) to equilibrium activities (an equal-affinity reference state).

The construction of 2-D predominance diagrams (for proteins or any other type of system) by default avoids calculating the equilibrium activities of species and instead identifies predominant species based on maximum affinity (after normalizing by the conservation coefficients). For systems of proteins, set `mam=FALSE` in `diagram()` to run the activity calculations if these values are needed, such as in the *E. coli* example above.

If oxygen fugacity is raised from its default nominal setting in CHNOSZ, the dynamic range of equilibrium activities calculated for proteins in human plasma becomes similar to the observed reference abundances of the proteins, and a slight positive correlation emerges. Equilibrium activities of kinases in *E. coli* cytosol have a dynamic range that is also similar to the observed abundances, but our findings so far imply going to a very low chemical potential of nitrogen (in terms of $\log a_{\text{NH}_3(aq)}$) to minimize the overall deviation.

7 Document Information

Revision history

- 2009-11-29 Initial version (Calculating relative abundances of proteins; current Section 2)
- 2011-06-20 Add human blood plasma and *E. coli* comparisons
- 2012-06-16 Use `protein.basis()` instead of `residue.info()`, `suppressMessages()` instead of `quiet=TRUE` in `info()`, `get.expr()` and `more.aa()`, and other minor updates.

R session information


```
> sessionInfo()

R version 2.15.0 (2012-03-30)
Platform: x86_64-slackware-linux-gnu (64-bit)

locale:
 [1] LC_CTYPE=en_US      LC_NUMERIC=C         LC_TIME=en_US
 [4] LC_COLLATE=C        LC_MONETARY=en_US    LC_MESSAGES=en_US
 [7] LC_PAPER=C          LC_NAME=C            LC_ADDRESS=C
[10] LC_TELEPHONE=C      LC_MEASUREMENT=en_US LC_IDENTIFICATION=C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  base

other attached packages:
[1] CHNOSZ_0.9-7.97

loaded via a namespace (and not attached):
[1] tools_2.15.0
```

References

- [1] N. L. Anderson and N. G. Anderson. The human plasma proteome - History, character, and diagnostic prospects. *Molecular & Cellular Proteomics*, 1(11):845–867, November 2002. doi: 10.1074/mcp.R200007-MCP200.
- [2] N. L. Anderson and N. G. Anderson. The human plasma proteome: History, character, and diagnostic prospects (vol 1, pg 845, 2002). *Molecular & Cellular Proteomics*, 2(1):50–50, January 2003. doi: 10.1074/mcp.A300001-MCP200.
- [3] J. M. Dick. Calculation of the relative metastabilities of proteins using the CHNOSZ software package. *Geochem. Trans.*, 9:10, 2008. doi: 10.1186/1467-4866-9-10.
- [4] J. M. Dick, D. E. LaRowe, and H. C. Helgeson. Temperature, pressure, and electrochemical constraints on protein speciation: Group additivity calculation of the standard molal thermodynamic properties of ionized unfolded proteins. *Biogeosciences*, 3(3):311 – 336, 2006. doi: 10.5194/bg-3-311-2006.
- [5] Jeffrey M. Dick and Everett L. Shock. Calculation of the relative chemical stabilities of proteins as a function of temperature and redox chemistry in a hot spring. *PLoS ONE*, 6(8):e22782, 2011. doi: 10.1371/journal.pone.0022782.
- [6] Y. Ishihama, T. Schmidt, J. Rappsilber, M. Mann, F. U. Hartl, M. J. Kerner, and D. Frishman. Protein abundance profiling of the *Escherichia coli* cytosol. *BMC Genomics*, 9:102, FEB 27 2008. ISSN 1471-2164. doi: 10.1186/1471-2164-9-102.
- [7] J. S. Seewald. Mineral redox buffers and the stability of organic compounds under hydrothermal conditions. *Mat. Res. Soc. Symp. Proc.*, 432:317 – 331, 1996. doi: 10.1557/PROC-432-317.