

Parameterisation of Food Balance Sheet Uncertainty Distribution

Michael. C. J. Kao

Food and Agriculture Organization
of the United Nations

Abstract

In solving the imbalance problem of the Food Balance Sheet (FBS), a novel solution based on probability maximisation was adopted to avoid the dispute of which element should be chosen to balance.

The implementation of the probabilistic framework requires consistent and logical specification of corresponding probability distribution for each of the corresponding in the Food Balance Sheet.

In this paper, we present the rationale and theory behind the framework. At the same time we also provide example based illustration to demonstrate the use of the **faoswsFlag** packages.

Keywords: Uncertainty, Distribution.

1. Introduction

In preparing the Food Balance Sheet (FBS), one of the most indispensable yet challenging operation is the balancing mechanism. Due to the imperfection of data collection, estimation and imputation in the real world, it is the norm that the FBS is unbalanced and does not satisfy the equality constraint at first sight. Thus, a balancing mechanism is essential for satisfying the equality between the demand and supply of the FBS.

In current practice, when imbalance exist, a variable is assigned the balancing item and the value of the balancing item is adjusted such that the equality constraint is satisfied. The choice of a variable as a balancing item often reflects the availability of data (or the lack of data), rather than a logical justification and empirical evidence. It is therefore not surprising that different SUA compilers/SUA approaches have chosen different variables as their balancing items. USDA's balances, for instance, use feed (and residual use) as the balancing item, while the FBS often use food to balance supply and demand. Conveniently, the XCBS approach often chooses whatever variable is not explicitly available. Clearly, none of these approaches renders a satisfactory solution to the problem and, no matter what variable is used as the balancing item, this variable is fraught with the measurement errors of all other variables. Given the fact that there is no a priori reason to assume that the measurement errors cancel out, the balancing item is bound to be the most inaccurate variable of the balance. Extending the logic to the Food Balance Sheets, using food as the balancing item would therefore be the least suitable solution.

Problems associated with the current approach motivated the research team to seek a method in which inequality occurring in the FBS to be allocated to various elements based on a sound

and logical reasoning rather than arbitrary allocation. One method to handle the problem is to balance the FBS based on a probabilistic basis. Each element is assumed to have a pre-determined level of uncertainty, and the allocation of the imbalance will depend on the uncertainty of each element. That is, the smaller the uncertainty we have with a particular element, the less we should apportion the imbalance or adjust that particular element. In the extreme case where we have perfect certainty about an element, then no imbalance should be added to the element nor adjusted.

In order to proceed with the probabilistic balancing mechanism, formulation of distribution for each of the elements in the FBS is necessary. The specification of the distributions undermines the validity of the balancing mechanism, and thus a consistent and logical construction of the distribution is crucial. The distributions should reflect the underlying uncertainty associated with each element while respecting the relationship amongst all elements. It is under these conditions, the optimal solution from the probabilistic framework is valid and justified. A framework guiding the specification of the distributions and the parameterisations is thus required and is the focus of this paper. The absence of such framework generates inconsistency and paradox, further the use of the term probability maximisation is a disguise for a procedure which does not bear any proper interpretation.

The paper is organised as follows. A simplified dataset is presented to familiarise the user with the problem. The subsequent section is then devoted to the rationale and theory behind the method for constructing the uncertainty distribution; a simple example with snippets will accompany the theory section to elaborate on the methodology. Finally, an application of the method for the balancing mechanism is presented and ends with discussion.

2. The Data

Before we delve into the theory and application, we will demonstrate some cases of the data for background. Further, these data will be later used to demonstrate the method and the package.

First of all, we can load up the package by prompting the following command.

```
## Load the library and the two example datasets
library(faoswsFlag)
data(vignetteFlagTable)
data(vignetteFBSTable)
```

The example flag table *vignetteFlagTable* is shown below.

flagObservationStatus	flagObservationWeights
	1.00
I	0.10
E	0.02

The column **flagObservationStatus** represents the flags associated with each observed value, while the **flagObservationWeight** represents the corresponding confidence associated with the source. Empty flag denotes official data, while "I" stands for imputed value and "E" are manual estimates. The weights like probability, should be between 0 and 1.

Shown below is the *vignetteFBSTable* dataset which is a simplified version of the Food Balance

Sheet for illustration in this paper. Each observed value is associated with a flag which indicate the source of the data. Under the proposing framework, the flag contains information about the uncertainty of the value and will be used to parameterise the distribution. There are more elements to the Food Balance Sheet, however, we have selected a handful of variables for illustrative purpose.

production	flag	import	flag	export	flag	food	flag	loss	flag
220	I	10		50		150	E	100	I

In the example data, all trade data both import and export are official. On the other hand losses and production are imputed by statistical methods, while food was estimated based on manual estimates.

3. The Methodology

In this section, we will describe the rationale and provide some background theory followed by an example at the end.

From the example, we can observe that for each element and item in the Food Balance Sheet (FBS), only a single value is observed. Since only a single observation is available, the use of Frequentism method can not be applied here. Rather, we have adopted the subjectivism interpretation of probability in order to come up with a solution.

To construct a probability distribution about the value, one first requires a chosen particular distribution, then parameterise the distribution according to a set of standards and rules.

The choice of the distribution should reflect knowledge and known constraints about the variable. The support, shape and properties of the distribution should be guided by the expertise of the officer. For example, production is strictly positive and thus distributions such as the Normal or the Cauchy should be eliminated from the set. Further, if extreme value are more likely then a Weibull distribution may be more preferable in comparison to the truncated normal distribution which has a higher kurtosis.

After the distribution has chosen, then the distribution need to be parameterised in order to complete the construction of the distribution. Logically, the parameter of the distribution should relate to the observed value and the specified confidence. Moreover, the mode of the constructed distribution should be made to be equivalent to the observed value. This task is in general simple, yet the conversion of the confidence level to the dispersion parameter of the distribution requires several more steps.

Here we propose a method in which the confidence can be converted to a measure of uncertainty we have about a particular value and ultimately lead to the parameterisation of distributions.

3.1. Quantifying Uncertainty

To measure a piece of information, one can use the formula of *self information* which is a measure of the information content and is defined as,

$$I = -\ln(P) \quad P \in [0,1] \quad (1)$$

Where P is the probability or the confidence about the accuracy of the value assigned to the observed value in the first place. The natural logarithm is adopted here, but logarithm of any base can be used. This is a measure of the uncertainty conditioned on the confidence we have provided about the observed value. The greater the I , the larger the associated uncertainty, that is, the lower the confidence we assign to the particular value, the higher the uncertainty. When the confidence is 0, or with 0% certainty, then I is infinite or infinite uncertainty; on the other hand, when the confidence is 1 then the value is known with certain.

The logarithm also ensures that the uncertainty is additive. Essentially, the sum of the uncertainty is the log of the products of the probabilities assigned to the values. That is, it is the log of the joint probability assuming independence.

The function enable us to convert the confidence about a single value to how much uncertainty is associated with the value.

3.2. Parameterise Distribution Given Uncertainty

Provided that we observe a single value, and at the same time our quantifying the uncertainty about a particular value; any chosen distribution can be parameterised accordingly to reflect the empirical evidence and knowledge about the value.

By setting the observed value to the expected value of the distribution (the expected value here refers to the value with the highest probability, that is, the mode) and the self information to the expected information or the differential entropy of the distribution, the parameters of the chosen distribution can then be obtained by solving the set of equations. Given the level of uncertainty associated with each element, then regardless of the choice of distributions, one can always parametrize the distribution where the uncertainty is held the same. This provides a consistent framework for specifying distributions in which the uncertainty for each element is consistent and relative amongst all elements.

That is, we parameterise the distribution given the following identities.

$$\begin{aligned} Mode(X) &= x \\ H(X) &= I \end{aligned} \quad (2)$$

Where X is the random variable and x is the observed value, I is the self-information or uncertainty computed according to formula 1, and H is the differential entropy of the chosen distribution.

The main reason to use the entropy rather than other dispersion parameters such as the absolute size of the standard deviation is because it is unit free and does not depends on the size of the value. If we were to impose uncertainty between two values, then the uncertainty associated with both value should be set respectively to the confidence given independent of the magnitude of the value. For example, if we we have observed 2000 tonnes of wheat production and 1000 tonnes of food while the confidence in the two value are identical, then the

balance should be 1500 tonnes of production and food. If we were to base the uncertainty on standard deviation or percentage of variation, then the larger value will have large standard deviation of variation based on percentage and thus the final value will be closer to 1000 even though we have equal confidence in both values.

Furthermore, both Normal and the truncated Normal distribution has a standard deviation parameter, however, setting the two distribution with identical standard deviation actually gives the normal distribution a high level of uncertainty.

3.3. A Simple Example

The following illustration provides an example of the method, along with codes to demonstrate the use of the package. In addition, we will demonstrate how this framework can provide consistent parameterization of various distribution while maintaining the same level of uncertainty with the value.

```
obsValue = 20
confidence = 0.02
```

Let us assume that we have an estimated value of 20 thousand tonnes of wheat production in Australia in 2010.

Then following the flag table, we have a confidence of 2% in the observed value. The amount of uncertainty regarding the wheat production in Australia given the confidence can then be calculated as,

$$I = -\ln(0.02) \approx 3.912023$$

```
(selfInfo = selfInformation(confidence))

## [1] 3.912023
```

In order to preserve this uncertainty associated with this piece of information, we need to preserve the entropy of the distributions. That is, regardless which distribution we choose we need to parameterise the distribution such that the entropy is equivalent to the same nat of information available.

Now for naive reasons that we want to impose a Normal distribution on the wheat production in Australia, we can first set the observed value to the mode of the distribution to first give us the first parameter of the Normal distribution.

$$\mu = 20 \tag{3}$$

In order to solve for the standard deviation σ of the Normal distribution, we first re-arrange the entropy function of the Normal distribution where the parameter σ is a function of the

entropy H . Then by substituting the expected information H with the self-information of the observed value I , we can then obtain the standard deviation of the distribution

Starting with the differential entropy of the Normal distribution

$$H = \frac{1}{2} \ln(2\pi e \sigma^2)$$

and re-arrange the equation,

$$\sigma = \sqrt{\frac{e^{2H}}{2\pi e}}$$

substituting H with I we obtain the value of the standard deviation as

$$\sigma = \sqrt{\frac{e^{2(-\ln(0.02))}}{2\pi e}} \approx 12.0985$$

or simply,

```
parameterise(obsValue = obsValue,
             selfInformation = selfInfo,
             distribution = "normal")

## $mean
## [1] 20
##
## $sd
## [1] 12.09854
```

That is, when the observed value of wheat production in Australia is 20 and a confidence of 2% is imposed, then the associated uncertainty distribution is then:

$$W \sim N(20, 12.0985)$$

However, if we believe that the production is in general rather stable over time but are subject to events such as drought that can create extreme values, then the Cauchy distribution may be a more reasonable distribution. The same method also allow us to parameterise the Cauchy distribution in which the uncertainty remains constant. The choice of distribution should reflect our belief in the probability allocation but it should not alter the uncertainty we have imposed initially.

Following the same procedure, we obtain the following parameters.

$$x_0 = 20$$

$$\gamma = e^{I - \ln(4\pi)} = e^{-\ln(0.02) - \ln(4\pi)} \approx 3.9789$$

then,

$$W \sim \text{Cauchy}(20, 3.9789)$$

```
parameterise(obsValue = obsValue,
             selfInformation = selfInfo,
             distribution = "cauchy")

## $location
## [1] 20
##
## $scale
## [1] 3.978874
```

Note, since the standard deviation of the Cauchy distribution is undefined and thus it is impossible to parameterise the distribution if we based our uncertainty measure on the size of the standard deviation.

Moreover, we know production can not be negative and thus distributions such as the Normal or the Cauchy distribution with unbounded support may not be the appropriate distribution. A truncated Normal distribution may incorporate this information by truncating support and allow the variable to be defined strictly on the positive real line.

Following the principle, we will arrive at a distribution which preserves the uncertainty yet re-assign the probability to reflect the physical condition that production can not be negative. In the case of the truncated Normal distribution, analytical solution does not exist, but a numerical solution is provided by the package.

```
parameterise(obsValue = obsValue,
             selfInformation = selfInfo,
             distribution = "truncNorm")

## $mean
## [1] 20
##
## $sd
## [1] 15.02336
```

and the resulting distribution is,

$$W \sim \text{trN}(20, 15.0234)$$

When the mean is close to zero, the standard deviation of the truncated normal is marginally larger than the normal distribution above. This is due to the fact to maintain the same level of uncertainty while reducing the support space, one has to increase the standard deviation.

However, as the mean increases, the truncated normal becomes more like the normal distribution with very similar standard deviation.

Finally, the log-Normal distribution is also another distribution which is defined only on the real line that is suitable to describe the probability allocation of the wheat production.

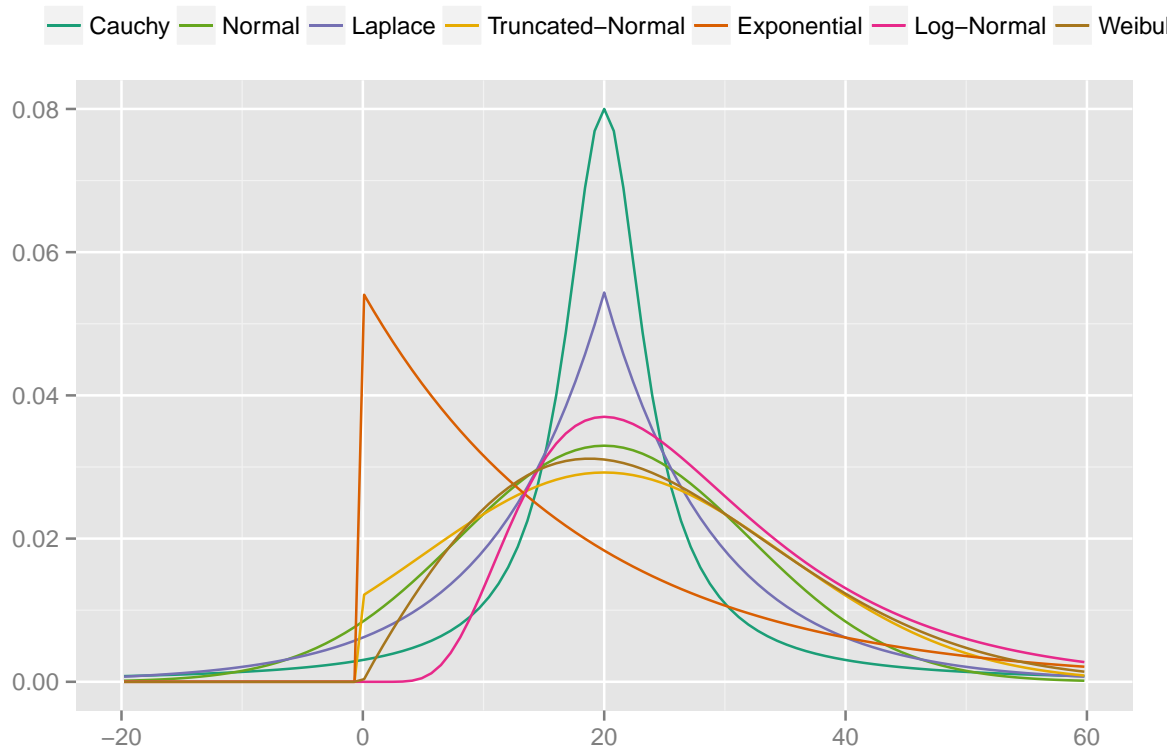
```
parameterise(obsValue = obsValue,
             selfInformation = selfInfo,
             distribution = "logNorm")
```

```
## $meanlog
## [1] 3.226412
##
## $sdlog
## [1] 0.4802915
```

and,

$$W \sim \ln N(3.2264, 0.4803)$$

Presented below is a comparison of the constructed distributions based on the same observed value and uncertainty. We can see all the distribution has mode at the observed value 20, except for the exponential distribution which has mode at 0. The choice of the distribution reflects our understanding of the element, but the level of uncertainty is held the same following the framework.



The log-Normal and the truncated-Normal distribution both converges to the Normal distribution when the mean is relatively large with respect to the standard deviation.

4. Illustration

To close the paper, we present a full case study of constructing the distributions and utilise the distributions to perform the probabilistic balancing.

Take the example data, the first step is to compute the level of uncertainty of each FBS element based on the flags in the FBS.

```
## Select all the flag columns
flagColumns = grep("Flag", colnames(vignetteFBSTable), value = TRUE)
valueColumn = grep("Value", colnames(vignetteFBSTable), value = TRUE)

## First we convert the flags to weights or confidence
(weightsFBS.df =
  data.frame(lapply(vignetteFBSTable[, flagColumns],
    function(x) {
      flag2weight(flagObservationStatus = x,
        flagTable = vignetteFlagTable)
    })))

##   productionFlag importFlag exportFlag foodFlag lossFlag
## 1             0.1           1           1      0.02      0.1

## Then we compute the self-information
(selfInfoFBS.df = data.frame(lapply(weightsFBS.df, selfInformation)))

##   productionFlag importFlag exportFlag foodFlag lossFlag
## 1      2.302585           0           0  3.912023  2.302585
```

To create the distribution, we simply provide the function `distributionise` the **observed value**, the **total information** computed from the flag and the **desired distribution**.

The function returns a list of two object. The first is the distribution function with the parameters computed, while the second object is a list with the corresponding values of the parameter.

```
## Parameterise the production element with a Normal distribution
distributionise(obsValue = vignetteFBSTable$productionValue,
  selfInformation = selfInfoFBS.df$productionFlag,
  distribution = "normal")

## $pdf
## function (x)
## dnorm(x, mean = mean, sd = sd)
## <environment: 0x5bfef48>
```

```
##
## $parameters
## $parameters$mean
## [1] 220
##
## $parameters$sd
## [1] 2.419707

## Parameterise the production element with a Truncated Normal distribution
distributionise(obsValue = vignetteFBSTable$productionValue,
               selfInformation = selfInfoFBS.df$productionFlag,
               distribution = "truncNorm")

## $pdf
## function (x)
## dtruncnorm(x, a = 0, b = Inf, mean = mean, sd = sd)
## <environment: 0x5a58368>
##
## $parameters
## $parameters$mean
## [1] 220
##
## $parameters$sd
## [1] 2.419707
```

Below we show a full process of how to construct each uncertainty distribution and specify the constraints for the balancing of the FBS.

```
## Here we simplify the example with one distribution,
## in practice each element can have their own corresponding distribution.
chosenDistribution = "truncNorm"

productionDist =
  distributionise(obsValue = vignetteFBSTable$productionValue,
                selfInformation = selfInfoFBS.df$productionFlag,
                distribution = chosenDistribution)

importDist =
  distributionise(obsValue = vignetteFBSTable$importValue,
                selfInformation = selfInfoFBS.df$importFlag,
                distribution = chosenDistribution)

exportDist =
  distributionise(obsValue = vignetteFBSTable$exportValue,
                selfInformation = selfInfoFBS.df$exportFlag,
                distribution = chosenDistribution)

foodDist =
  distributionise(obsValue = vignetteFBSTable$foodValue,
```

```

        selfInformation = selfInfoFBS.df$foodFlag,
        distribution = chosenDistribution)
lossDist =
  distributionise(obsValue = vignetteFBSTable$lossValue,
    selfInformation = selfInfoFBS.df$lossFlag,
    distribution = chosenDistribution)

## Create the likelihood function from the distributions
ll = function(x){
  -log(productionDist$pdf(x[1])) -
    log(importDist$pdf(x[2])) -
    log(exportDist$pdf(x[3])) -
    log(foodDist$pdf(x[4])) -
    log(lossDist$pdf(x[5]))
}

## Create the constraint function
constraint = function(x){
  ## (1) Production + Import - Export - Food - Loss = 0
  ## (2) and (3) holding import and export constant
  c(x[1] + x[2] - x[3] - x[4] - x[5], x[2], x[3])
}

## Balance the Food Balance Sheet
library(Rsolnp)
## NOTE (Michael): If degenerate distribution is present, then take
##                  a small step (small delta) and have large inner
##                  iteration to ensure convergence to the right solution.
balancedFBS =
  solnp(pars = as.numeric(vignetteFBSTable[, valueColumn]),
    fun = ll,
    eqfun = constraint,
    eqB = c(0, vignetteFBSTable$importValue, vignetteFBSTable$exportValue),
    LB = rep(0, length(valueColumn)),
    control = list(delta = 1e-1, trace = 2, inner.iter = 3000))

##
## Iter: 1 fn: 22.5156 Pars: 222.60470 10.00000 50.00000 85.25698 97.34772
## Iter: 2 fn: 22.5156 Pars: 222.60404 10.00000 50.00000 85.25621 97.34784
## Iter: 3 fn: 22.5156 Pars: 222.60398 10.00000 50.00000 85.25613 97.34785
## Iter: 4 fn: 22.5156 Pars: 222.60392 10.00000 50.00000 85.25606 97.34786
## Iter: 5 fn: 22.5156 Pars: 222.60389 10.00000 50.00000 85.25603 97.34786
## solnp--> Completed in 5 iterations

## Original values
as.numeric(vignetteFBSTable[, valueColumn])

## [1] 220 10 50 150 100

```

```
## Balanced values  
round(balancedFBS$par)  
  
## [1] 223 10 50 85 97
```

From the balanced value, we can see since the element Food has the lowest confidence, thus its value is adjusted the most followed by Production and Loss with Import and Export held constant.

In addition, since we have the same confidence in both the production and the losses element while at the same time specified symmetric distributions; the pair were adjusted towards the balancing by the same magnitude.

Affiliation:

Michael. C. J. Kao

Economics and Social Statistics Division (ESS)

Economic and Social Development Department (ES)

Food and Agriculture Organization of the United Nations (FAO)

Viale delle Terme di Caracalla 00153 Rome, Italy

E-mail: michael.kao@fao.org

URL: https://github.com/mkao006/sws_r_api/tree/master/faoswsFlag