# CODEX ALIMENTARIUS COMMISSION

**Food and Agriculture Organization of the United Nations**

**World Health Organization**

**JOINT FAO/WHO FOOD STANDARDS PROGRAMME**

**CODEX COMMITTEE ON CONTAMINANTS IN FOODS**

**13th Session**

**Yogyakarta, Indonesia, 29 April – 3 May 2019**

**GENERAL GUIDANCE ON DATA ANALYSIS FOR ML DEVELOPMENT**

(*Prepared by EU*)

## BACKGROUND

1. At its 12th session, CCCF considered the proposal of the JECFA Secretariat to develop a general guidance on data analysis for ML development as it was observed that different approaches were taken by the EWGs. These differences concerned for example the handling of occurrence data without information on LOQ. A general guidance would help future EWGs to take consistent approaches for data analysis. CCCF agreed to establish an EWG chaired by EU, co-chaired by the United States of America, the Netherlands and Japan, working in English, to prepare a discussion paper[1]

2. It has not been possible to prepare in time a discussion paper for consideration by the established EWG. Therefore, this discussion paper contains a non-exhaustive list of items, prepared by the chair of the EWG, that could be considered to be covered by the general guidance on data analysis for ML development. Based on a preliminary discussion at the current session, a more elaborate document outlining a general guidance on data analysis for ML development should be prepared for consideration by the EWG in view of a discussion at CCCF14.

## NON-EXHAUSTIVE LIST OF TOPICS POSSIBLY TO BE CONSIDERED FOR GENERAL GUIDANCE ON DATA ANALYSIS FOR ML DEVELOPMENT.

### 1) Criteria for the establishment of maximum levels in food and feed[2]

**Selection of criteria has been made of relevance for occurrence data as basis for setting MLs**

- Validated qualitative and quantitative analytical data on representative samples should be supplied. Information on the analytical and sampling methods used and on the validation of the results is desirable. A statement on the representativeness of the samples for the contamination of the product in general (e.g. on a national basis) should be added. The portion of the commodity that was analyzed and to which the contaminant content is related should be clearly stated and preferably should be equivalent to the definition of the commodity for this purpose or to existing related contaminant regulation.

- Information on appropriate sampling procedures should be supplied. Special attention to this aspect is necessary in the case of contaminants that may not be homogeneously distributed in the product (e.g. mycotoxins in some commodities).

- MLs should be set as low as reasonably achievable and at levels necessary to protect the consumer. Providing it is acceptable from the toxicological point of view, MLs should be set at a level which is (slightly) higher than the normal range of variation in levels in food and feed that are produced with current adequate technological methods, in order to avoid undue disruptions of food and feed production and trade. Where possible, MLs should be based on GMP and/or GAP considerations in which the health concerns have been incorporated as a guiding principle to achieve contaminant levels as low as reasonably achievable and necessary to protect the consumer. Foods that are evidently contaminated by local situations or processing conditions that can be avoided by reasonably achievable means shall be excluded in this evaluation, unless a higher ML can be shown to be acceptable from a public health point of view and significant economic aspects are at stake.

---

[1] REP18/CF, paras 155-156

[2] Reference is made to the criteria for the establishment of maximum levels in food and feed as provided for in Annex I of CXS 193-1995 General Standard for Contaminants and Toxins in Food and Feed

- Proposals for MLs in products should be based on data from various countries and sources, encompassing the main production areas/processes of those products, as far as they are engaged in international trade. When there is evidence that contamination patterns are sufficiently understood and will be comparable on a global scale, more limited data may be enough.

- MLs may be set for product groups when sufficient information is available about the contamination pattern for the whole group, or when there are other arguments that extrapolation is appropriate.

- Numerical values for MLs should preferably be regular figures in a geometric scale (0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1, 2, 5 etc.), unless this may pose problems in the acceptability of the MLs.

- MLs should apply to representative samples per lot. If necessary, appropriate methods of sampling should be specified.

## 2) Minimum number of samples needed for the use of percentiles

### Background

In order to apply the above criterion "*MLs should be set at a level which is (slightly) higher than the normal range of variation in levels in food and feed*", high percentiles are used to define that level. The reliability of high percentiles is related to the number of data used to calculate them. Percentiles calculated on a number of subjects should be treated with caution as the results may not be statistically robust.

A clear indication concerning the minimum number of observations necessary to estimate a given percentile is not provided in literature. Different options can be used, none of them being a widely accepted standard.

A very simple option is to require that the calculated percentile must at least be different from the maximum value within the sample. This means that at least 20 observations are needed to identify the single observation at the 95th percentile and 100 observations are needed for the 99th percentile.

In statistics, the coverage probability of a confidence interval is the probability that the interval contains the true value of interest (e.g. 95th or 99th percentiles). When the number of observations is not large enough, the coverage probability may not attain the nominal value, and drops below, for example, 95%. This is more likely to occur at high percentiles, e.g. 95th or 99th. Therefore, the coverage probability has been used to set guidelines to determine the minimum number of samples for which (extreme) percentiles can be computed. In the case of significance level ($\alpha$) being set at 0.05 to determine a 95% confidence interval, the coverage probability should target 95%. In this case, this is achieved for n ≥ 59 and n ≥ 298 for the 95th or 99th percentiles, respectively.

## 3) Limit of Quantification (LOQ) considerations

Several situations applicable to datasets provided can occur and the guidelines to be elaborated should provide guidance on how to handle the datasets in the different situations

- No LOQ provided

  o Dataset contains (nearly) all quantified results

  o Dataset contains a significant part of left-censored data (i.e. < LOQ) and no LOQ provided

In the above situations where the LOQ is not provided, should the guidance provide for different conclusions as regards how to handle the dataset in case the quantified results (significantly lower than the ML under consideration) in the dataset provide a an indication that the LOQ is (very) low compared to datasets where the quantified results do not provide that indication.

- LOQ provided

  o Dataset with LOQ significantly lower than the ML under consideration

  o Dataset with LOQ in the range of the ML under consideration

  o Dataset with LOQ above the ML under consideration

In the above situations where the LOQ is provided, should there be guidance on cut-offs to be used for the LOQ on the analytical results dataset used for the ML development?

Should the guidance provide for different conclusions as regards how to handle the dataset in case the dataset contains nearly all quantified results compared to a dataset with nearly all left-censored data?

**4) Using data sets with a large proportion of left-censored data for ML development**

In certain cases, the analytical results for one specific contaminant are produced with a battery of different analytical methods and/or the same analytical method but with very different sensitivities. As a consequence, there could be a wide range of limits of detection (LOD) and limits of quantification (LOQ) for a particular contaminant and food matrix in a given dataset, composed of datasets from different sources. This situation is particularly relevant when the occurrence datasets used for the ML development contain a high number of non-quantified/non-detected data (left-censored data).

The standard approach to deal with left-censored data is the use of the substitution. In this method, at the lower-bound (LB), results below the LOQ and LOD are replaced by zero; at the upper-bound (UB) the results below the LOD are replaced by the numerical value of the LOD and those below the LOQ are replaced by the value reported as LOQ. Additionally, as a point estimate between the two extremes, the middle-bound (MB) scenario is calculated by assigning a value of LOD/2 or LOQ/2 to the left-censored data.

**5) Geographical coverage of the provided occurrence data**

Guidance should be provided to evaluate the appropriateness of the geographical coverage of the provided data for ML development and a procedure should be developed for situations for which it is concluded that the available data do not provide a sufficient/appropriate geographical coverage;

**6) Period coverage of the provided occurrence data**

Guidance should be provided in which situation it might be required that the provided occurrence data relate to several production years for ML development (can be different for different types of contaminants: mycotoxins, plant toxins, processing contaminants, environmental contaminants in function of the assumed year to-year variation or evolution of contamination in time)

**7) Evaluation if provided occurrence data reflect the application of Codex Code of Practice and/or GAP/GMP**

Consideration could be given if it is possible to provide guidance on which elements a possible evaluation should be based to determine if provided datasets do reflect the application of a Codex Code of Practice and/or GAP/GMP.

**8) Data sets with low number of data (e.g. less than 60) for development of ML**

Guidance could be given in which situations it can be concluded that the data, despite the low number, are sufficient for the development of an ML (e.g. despite limited number good geographical coverage, no large variation in occurrence observed despite data originating from different regions/from different years, etc).

**9) Other elements of data analysis for ML development for which guidance should be developed.**

**RECOMMENDATIONS**

3.  CCCF13 is invited to

    a)  Have a preliminary consideration which of the abovementioned identified topics are appropriate for inclusion in a guidance for data analysis for ML development

    b)  Have a non-binding consideration of other topics which would be appropriate to be included in a guidance for data analysis for ML development

    c)  Agree to re-establish the EWG to prepare a draft of a general guidance on data analysis for ML development taking into account the outcome of the preliminary discussion at this meeting fro discussion at CCCF14 (2020).