

# Meeting Summary for the Live Professor Q&A on Model Fitting, Bootstrap, and Model Selection at the 2025 Astrostatistics Summer School (06/03/2025)

## 1 Quick recap

The meeting began with technical difficulties as Jogesh attempted to share his screen to explain confidence bands, followed by G.'s detailed explanation of computing confidence intervals and the empirical distribution function. The discussion then covered various statistical concepts including boosting methods, the Law of Iterated Logarithm, and spectral line modeling, with participants asking clarifying questions throughout. The session concluded with an exploration of parametric bootstrap methods for estimating distribution functions and calculating the Kolmogorov-Smirnov test statistic, particularly in cases of unknown parameters and heteroscedastic data.

## 2 Summary

### 2.1 Screen Sharing Technical Difficulties

Jogesh attempted to share his screen to explain Ks confidence bands but encountered technical difficulties due to switching between Wi-Fi networks. He asked the participants if they could see his screen and mentioned he was having trouble with audio.

### 2.2 Confidence Interval Calculation Process

G. explained the process of computing a confidence interval using the  $D_n$  Alpha statistic, which involves finding a value  $Y$  such that the probability of the  $D_n$  statistic being less than or equal to  $Y$  is equal to the desired confidence level (e.g., 0.95). This value is then added to and subtracted from the empirical distribution function to create a confidence band. Sagnik asked for clarification, and G. repeated the explanation, emphasizing the steps to find  $D_n$  Alpha and how to use it to construct the confidence interval.

### 2.3 Empirical Distribution Function Overview

G. explained the concept of the empirical distribution function, which is a step function that represents the cumulative distribution of data. Jay asked for clarification on how to calculate the empirical distribution function, and G. described the process of counting

observations less than or equal to a given value and dividing by the total number of observations. G. also explained that this function is similar to the cumulative distribution function (CDF) and is used to normalize data into a probability distribution.

## **2.4 X vs X Star Calculation**

Shlok asked about the difference between X and X Star, and G explained that X Star involves multiplying the number on a chip by the number of chips, while X involves multiplying the number on a chip by 1. G clarified that the number of elements and values are the same in both cases, and there may be repetitions in the X Star process. G also mentioned that they would address questions about the "booster fail" slide later.

## **2.5 Limitations of Boosting in Heavy-Tails**

G. explained that boosting methods work well for most statistical analyses involving smooth data, but they fail in cases of heavy-tailed distributions and infinite variance. G. attributed this limitation to the concept of adverse expansions, which is too complex to explain at this level. They also noted that the maximum value in heavy-tailed distributions does not change smoothly with new observations, making it difficult for boosting methods to match the distribution accurately.

## **2.6 Bootstrap Analysis and Convergence Laws**

G. explained the practical use of the Law of Iterated Logarithm in analyzing bootstrap applications, which involve  $N \log N$  squared computations. Sagnik sought clarification on how the number of applications suffices, which G. attributed to the Law of Iterated Logarithms and the fast convergence of error terms in bootstrap distributions. They also discussed the distinction between continuous and discrete distributions, emphasizing that even with large sample sizes, discrete distributions like binomial or Poisson remain discrete.

## **2.7 Spectral Line Modeling and Analysis**

The discussion focused on the properties of spectral lines and their mathematical modeling. G. explained that while natural line shapes are typically Lorentzian, thermal Doppler broadening dominates observed line widths, and the Y function model is a convolution of Cauchy, Lorentzian, and Gaussian functions. They noted that bootstrap

approaches may not be effective for inferring line parameters when dealing with Cauchy distributions, as these do not follow the central limit theorem. G. also mentioned that parametric and nonparametric bootstraps provide similar results in infinite-dimensional cases like Kolmogorov's work.

## **2.8 Parametric Bootstrap for Heteroscedastic Data**

The discussion focused on the use of parametric bootstrap methods for estimating distribution functions and calculating the Kolmogorov-Smirnov test statistic, particularly when dealing with unknown parameters and heteroscedastic data. It was explained that simple bootstrap methods are not suitable for heteroscedastic data, and more sophisticated approaches are needed. The speaker emphasized that the bootstrap method essentially calculates the Kolmogorov-Smirnov table and recommended using studentized bootstrap for normalization when the variance is unknown. They encouraged participants to ask further questions via email or the Slack channel.