## *14.3 Are Two Distributions Different?*

Given two sets of data, we can generalize the questions asked in the previous section and ask the single question: Are the two sets drawn from the same distribution function, or from different distribution functions? Equir fromconsistentwith a single distributionw( [(distrib)24.1(uted)-301.1(uniform brands, or Brooklyn and the Bronx).

One can always turn continuous data into binned data, by grouping the events into specified ranges of the continuous variable(s): declinations between 0 and 10 degrees, Binningvolves a loss of however.

Also, there is often considerable arbitrariness as to how the bins should be chosen. Along with many other investigators, we prefer to avoid unnecessary binning of data.

The accepted test for differences between binned distributions is the *chi-square test*. For continuous data as a function of a6(a)-409single variable, most accepted test is the *Kolmogorov-Smirnov testKe*xoneptedisconeptedvestigatoisispre81

integers, while the $n_i$'s may not be. Then the chi-square statistic is

$$\chi^2 = \sum_i \frac{(N_i}{}$$