

1116-62-1317

Xiao-Li Meng* (meng@stat.harvard.edu). *Statistical Paradises and Paradoxes in Big Data*.

Statisticians are increasingly posed with thought-provoking and often paradoxical questions, challenging our qualifications for entering the statistical paradises created by Big Data. Questions addressed in this talk include 1) Which one should I trust: a 1% survey with 60% response rate or a self-reported administrative dataset covering 80% of the population? 2) With all the big data, is sampling or randomization still relevant? 3) Personalized treatments—that sounds heavenly, but where on earth did they find the right guinea pig for me? The proper responses are respectively 1) “It depends!” because we need *data-quality indexes*, not merely quantitative sizes, to determine; 2) “Absolutely!” and indeed Big Data has inspired methods such as *counterbalancing sampling* to combat inherent selection bias in big data; and 3) “They didn’t!” but the question has led to a *multi-resolution framework* for studying statistical evidence for predicting individual outcomes. All proposals highlight the need, as we get deeper into this era of Big Data, to reaffirm some time-honored statistical themes (e.g., bias-variance trade-off), and to remodel some others (e.g., approximating individuals from proxy populations verses inferring populations from samples). (Received September 18, 2015)