

The rise of massive datasets that provide fine-grained information about human beings and their behavior provides unprecedented opportunities for evaluating the effectiveness of social, behavioral, and medical treatments. We have also become more interested in fine-grained inferences. Researchers and policy makers are increasingly unsatisfied with estimates of average treatment effects based on experimental samples that are unrepresentative of populations of interest. Instead, they seek to target treatments to particular populations and subgroups. These fine-grained inferences lead to small data problems: subgroups where the dimensionality of data is high but the number of observations is small. To make the best use of these new methods, randomized trials should be constructed differently, with an eye towards how they will be combined with observational data down the road. I discuss new methods for designing and analyzing randomized experiments that make the most of these opportunities. For example, inferences from randomized experiments can be improved by blocking: assigning treatment in fixed proportions within groups of similar units. However, the use of the method is limited by the difficulty in deriving these groups. Current blocking methods are restricted to special cases or run in exponential time; are not sensitive to clustering of data points; and are often heuristic, providing an unsatisfactory solution in many common instances. We present an algorithm that implements a new, widely applicable class of blocking—threshold blocking—that solves these problems. Given a minimum required group size and a distance metric, we study the blocking problem of minimizing the maximum distance between any two units within the same group. We prove this is a NP-hard problem and derive an approximation algorithm that yields a blocking where the maximum distance is guaranteed to be at most four times the optimal value. This algorithm runs in $O(n \log n)$ time with $O(n)$ space complexity. This makes it the first blocking method with an ensured level of performance that works in massive experiments. While many commonly used algorithms form pairs of units,

our algorithm constructs the groups flexibly for any chosen minimum size. This facilitates complex experiments with several treatment arms and clustered data. I also discuss extensions of this method for observational data (e.g., matching) and for exploratory data analysis (e.g., clustering).