

# Indiana University-Purdue University Indianapolis

Department of Mathematical Sciences

STATISTICS SEMINAR

12:15pm—1:15pm, Tuesday, Feb. 13, 2018

SL 137

**Speaker:** **Wei Zhao** (PhD Candidate)  
*Department of Mathematical Sciences, IUPUI*

**Title:** **Information-Based Optimal Subdata Selection for Big Data Linear Regression**

## **Abstract:**

Extraordinary amounts of data are being produced in many branches of science. Proven statistical methods are no longer applicable with extraordinary large data sets due to computational limitations. A critical step in big data analysis is data reduction. Existing investigations in the context of linear regression focus on subsampling-based methods. However, not only is this approach prone to sampling errors, it also leads to a covariance matrix of the estimators that is typically bounded from below by a term that is of the order of the inverse of the subdata size. We propose a novel approach, termed information-based optimal subdata selection (IBOSS). Compared to leading existing subdata methods, the IBOSS approach has the following advantages: (i) it is significantly faster; (ii) it is suitable for distributed parallel computing; (iii) the variances of the slope parameter estimators converge to 0 as the full data size increases even if the subdata size is fixed, i.e., the convergence rate depends on the full data size; (iv) data analysis for IBOSS subdata is straightforward and the sampling distribution of an IBOSS estimator is easy to assess. Theoretical results and extensive simulations demonstrate that the IBOSS approach is superior to subsampling-based methods, sometimes by orders of magnitude.

## **Reference:**

Wang, HaiYing, Min Yang, and John Stufken. "Information-Based Optimal Subdata Selection for Big Data Linear Regression." *Journal of the American Statistical Association* just-accepted (2017).