

Bag of Little Random Forests (BLRF)

Adapting the RF to a Big Data Setting

Motivation and Objectives

The Random Forests (RF) algorithm (Figure 1) is inefficient when handling Big Data:

- ❖ **Time**: significant time consumption;
- ❖ **Memory**: physical storage of big data sets;
- ❖ **Structure**: not well-adapted to a parallelism;
- ❖ **Unable** to load the entire data set into memory.

Aim to: build BLRF (Figure 3). Reduce computation time while maintaining prediction accuracy.

Steps and Procedures

- ❖ Study BLB and RF;
- ❖ Combine BLB and RF -> BLRF algorithm;
- ❖ Modify the source C code (regression only);
- ❖ Evaluate the BLRF algorithm (time and accuracy).

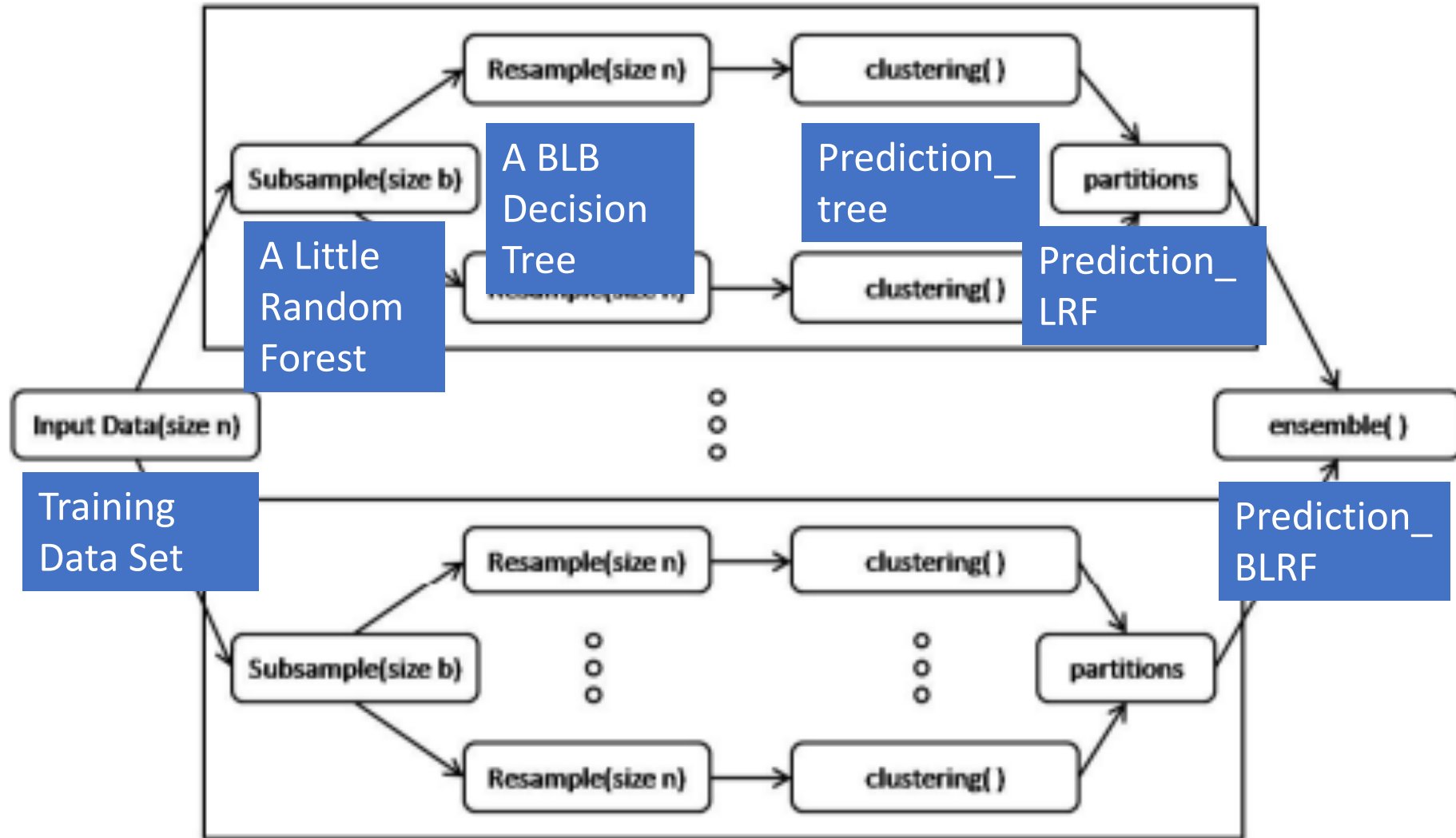


Figure 3. Visualization of BLRF

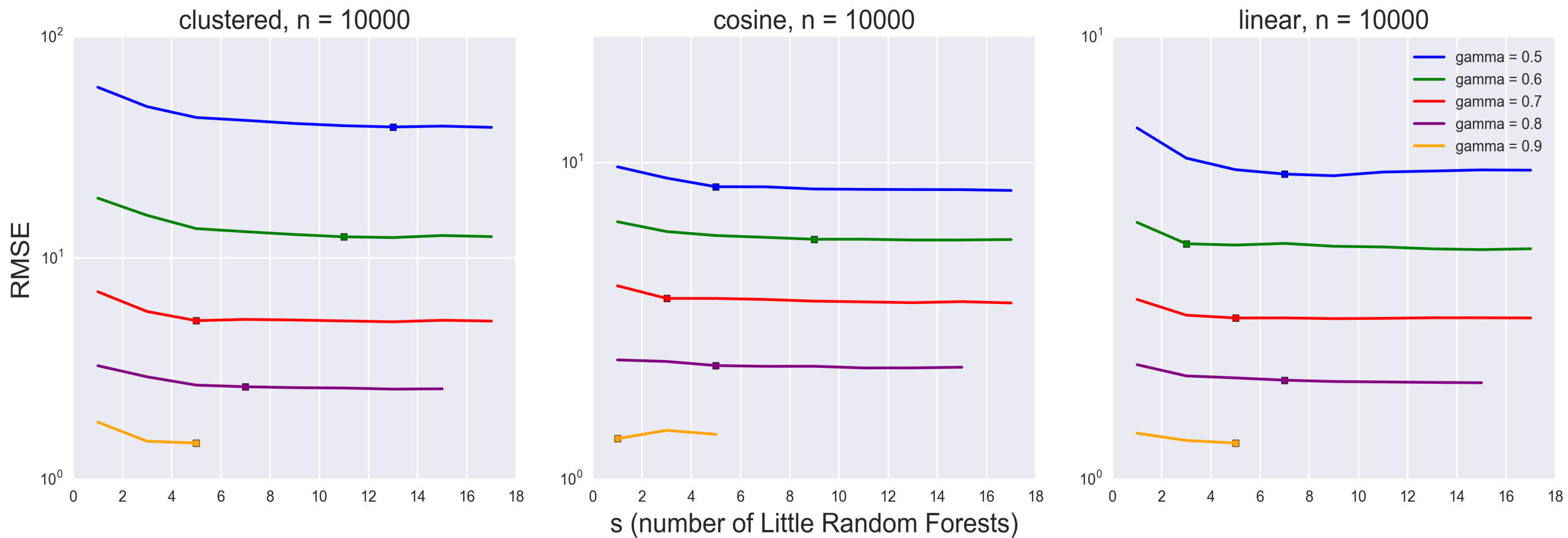
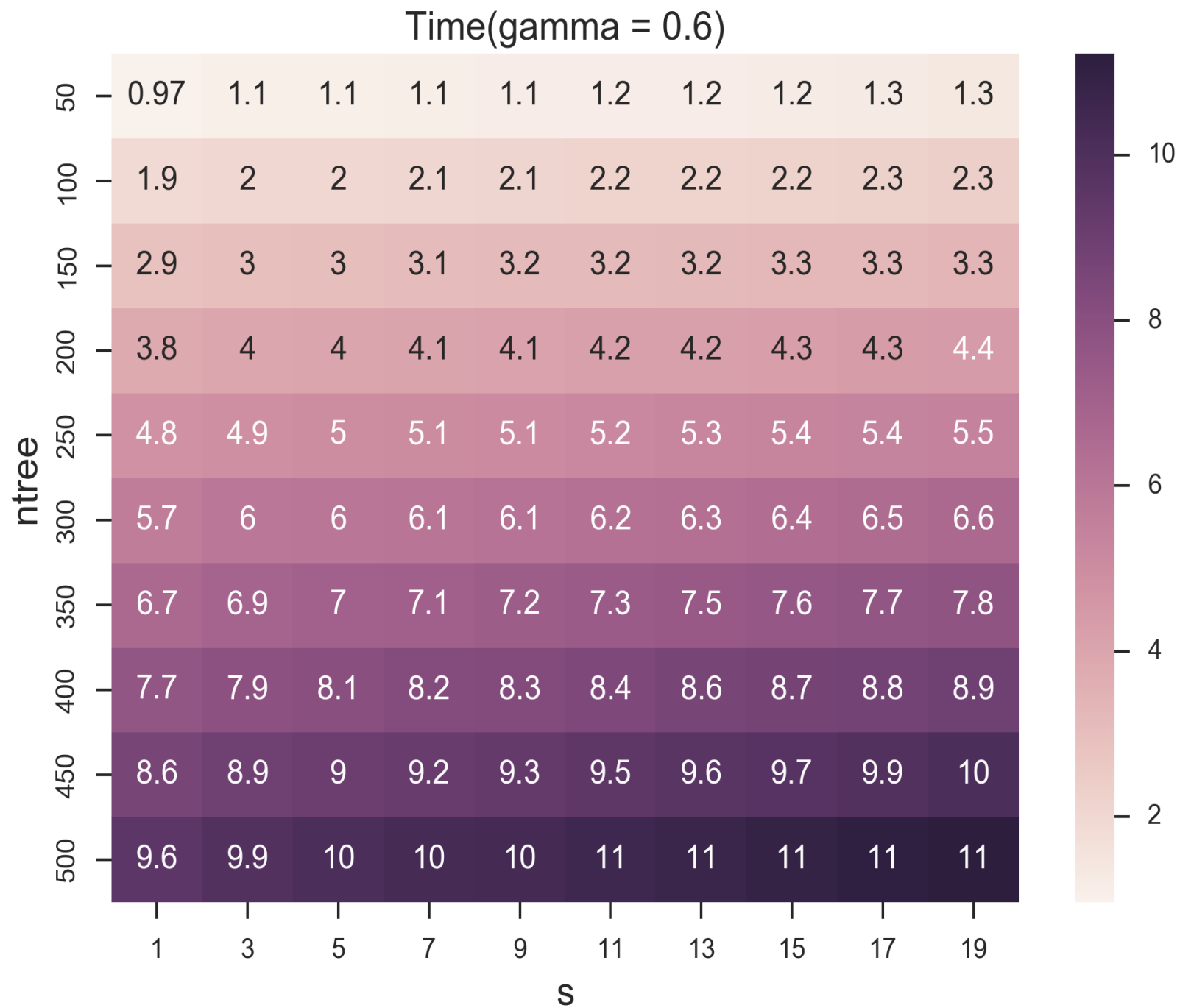


Figure 5. RMSE measures of BLRF algorithm: $\text{RMSE} \sim s$



Performance of BLB-RF vs RF

