BUAN 6346 Big Data

## Homework 10b

# Explore RDDs using the Spark Shell

**Step11:** mydata.collect()