

Some Background

Lasso regression:

- Regularization is used to avoid **overfitting**.
- Lasso: L_1 penalty leads to implicit feature selection.
- Tuning of λ with cross-validation.
- Comparison: *glmnet* R-package vs. *LassoWithSGD*.

Algorithm implemented in *glmnet*: Friedman, Jerome, Trevor Hastie, Holger Höfling, and Robert Tibshirani. "Pathwise coordinate optimization." The Annals of Applied Statistics 1, no. 2 (2007): 302-332.

Apache Spark:

- A general purpose engine for large-scale data processing.
- MLlib**: distributed machine learning algorithms.
- Can process data from local file system, HDFS, distributed data bases, or streams.



Methodology

Instance generator:

- Number of observations N .
- Number of features d .
- Ratio $\frac{N}{d}$.
- Signal-to-noise ratio $n \in [0, 1]$.
- Polynomial degree p .
- Generative mechanism for $p = 2$:
 $y = N(c_{x11}x_1 + c_{x12}x_1^2 + \dots + c_{xd1}x_d + c_{xd2}x_d^2, 1 - n)$

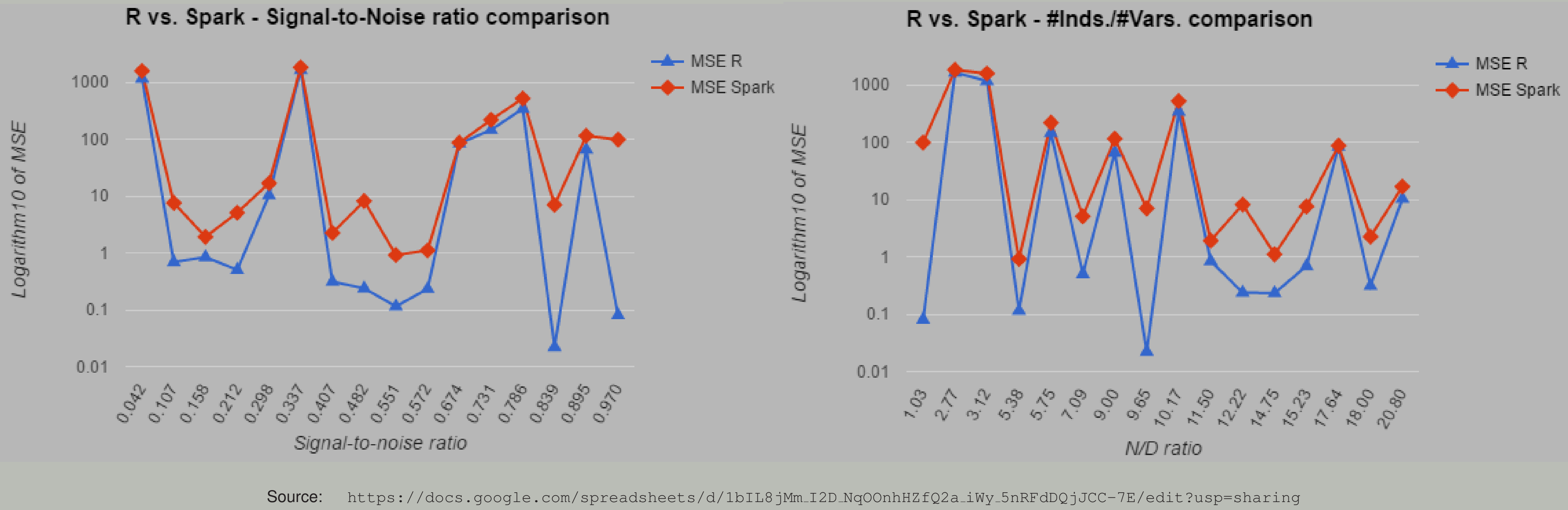
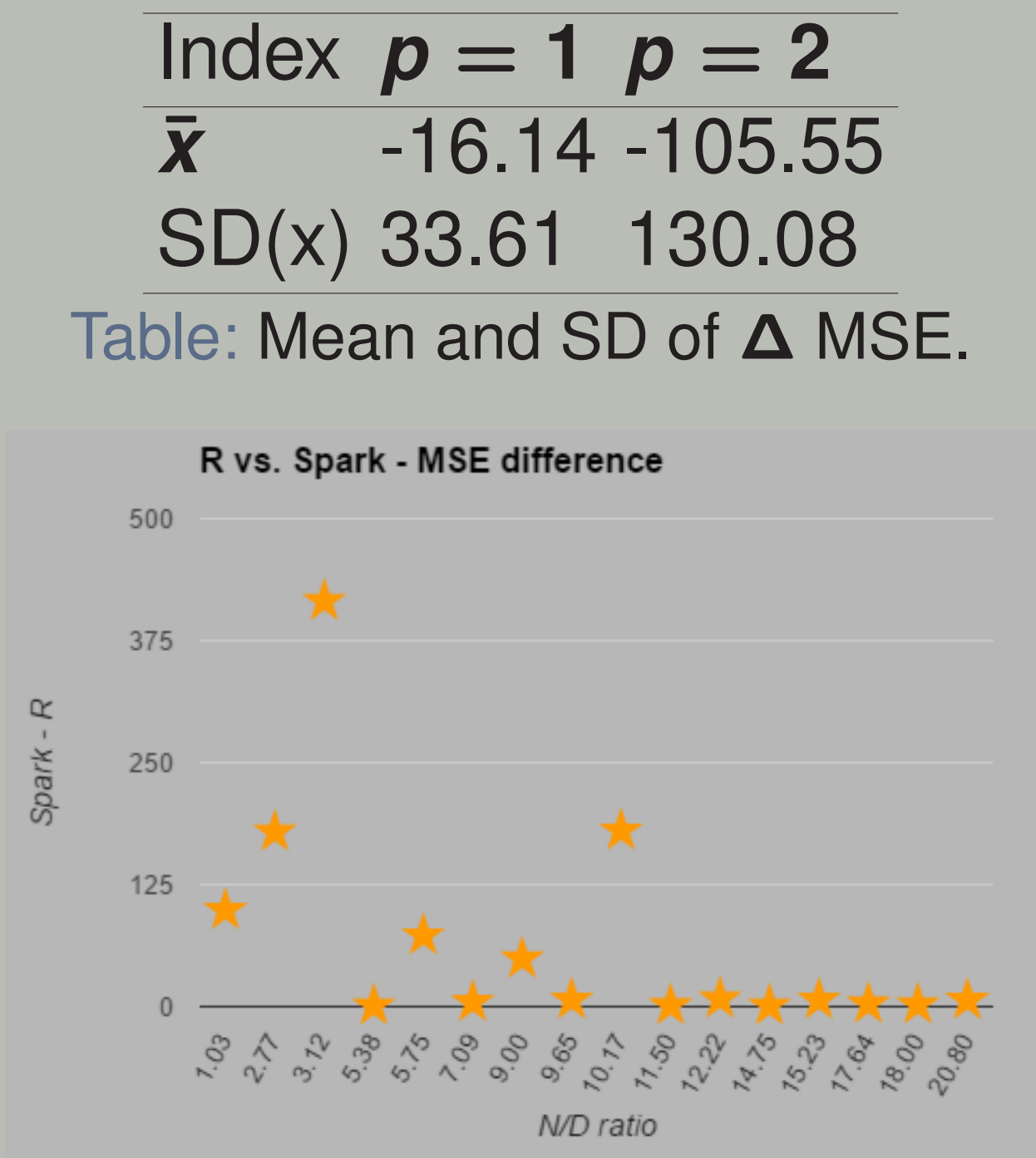
Latin Hypercube sampling:

- n : signal-to-noise ratio,
- N : number of observations,
- $\frac{N}{d}$: observations / features,
- p : polynomial degree.

Methodological objective:

- Quantify complexity of data sets.

Results



Motivation

There is a hype about “big data analytics”:

- Hadoop** (Mahout),
- Spark** (MLlib),
- Flink** (FlinkML).

But, how **accurate** and **performant** are machine learning algorithms based on distributed programming models? How does their **usability** compare with standard, centralised algorithms available e.g. in *R*?

Motivation

In this project we compared the Lasso (L1) regression model between R and MLlib (Spark) in terms of:

- accuracy,
- performance,
- usability.

We used 16 synthetic data sets differing in difficulty.

Conclusions: Usability

Compared to *R*, in Spark MLlib:

- limited range of machine learning methods.
- only token implementations of some methods (i.e. SVM).
- manipulating data interactively in shell is clumsy.
- parameter-tuning requires sophisticated code.
- debugging is more difficult.

Conclusions: Accuracy and Performance

- The *glmnet* algorithm computes optimal λ in one go!
- MSE of MLlib was **always** higher for a given data set!
- The error of MLlib (Spark) increased with $\frac{N}{d}$ and p .
- MLlib can not deal with difficult data sets!

Acknowledgment

This poster was reviewed by Santiago Rodrigo Muñoz and Hao Wu.