

Least squares method for Time series

Compare 3 methods:

- Pseudo inverse
- QR
- SVD

Problem description

$$\hat{x}_{k+1} = w_1 x_k + \dots + w_M x_{k-M+1}$$

$$\|Xw - y\|_2^2 \rightarrow \min_w$$

The goal is to compare 3 methods: pseudo inverse, QR and SVD

Toeplitz

- [Toeplitz matrix - Wikipedia](#)
- $O(n)$ add
- $O(n \log n)$ vector mult
- $O(n^2)$ matrix mult
- $O(n^2)$ to solve SLAE
- $O(n^2)$ decomposition
- scipy works only with square matrices :(

$$A = \begin{bmatrix} a_0 & a_{-1} & a_{-2} & \cdots & \cdots & a_{-(n-1)} \\ a_1 & a_0 & a_{-1} & \ddots & & \vdots \\ a_2 & a_1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & a_{-1} & a_{-2} \\ \vdots & & \ddots & a_1 & a_0 & a_{-1} \\ a_{n-1} & \cdots & \cdots & a_2 & a_1 & a_0 \end{bmatrix}$$

Method of comparison of algorithms

1. Splitting a dataframe with different ratios
2. Training on a training samples
3. Calculating MSE on a test samples
4. Comparison MSE

Predict driver_pay

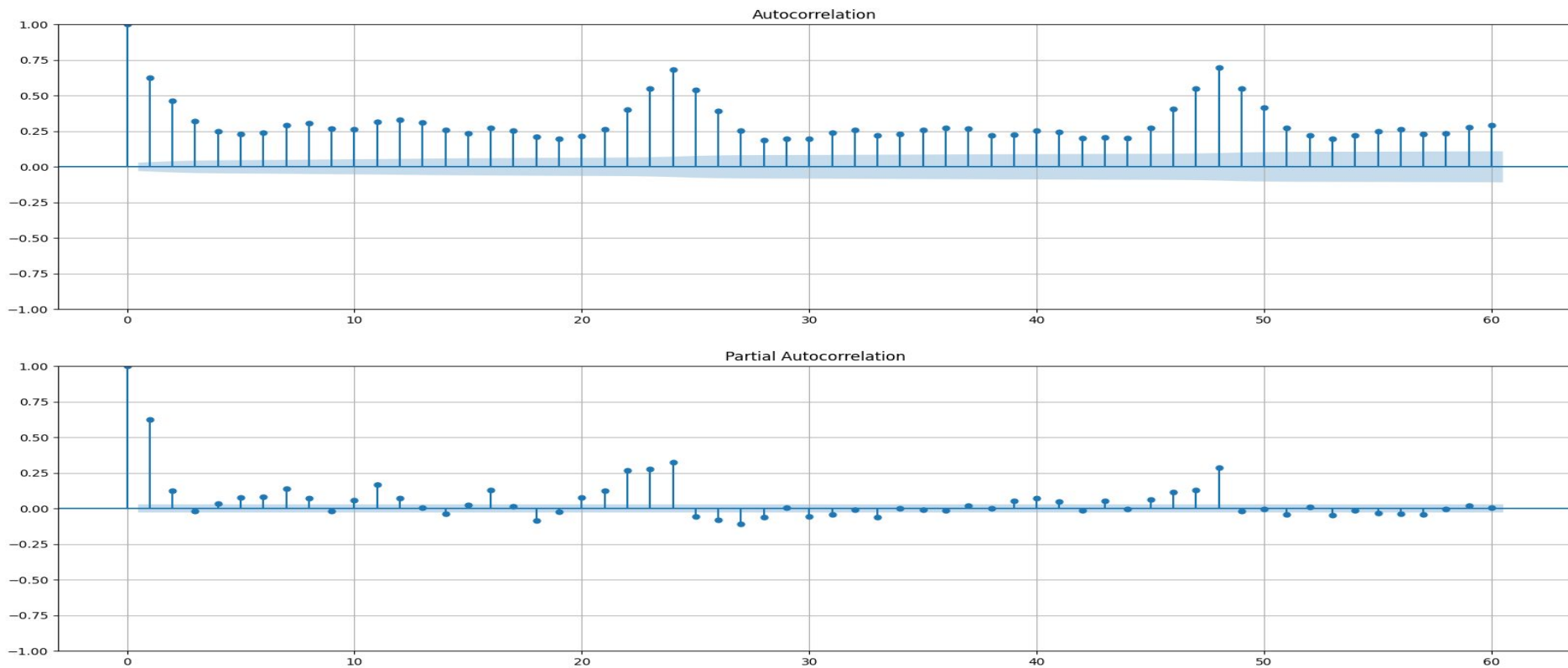
<https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

January

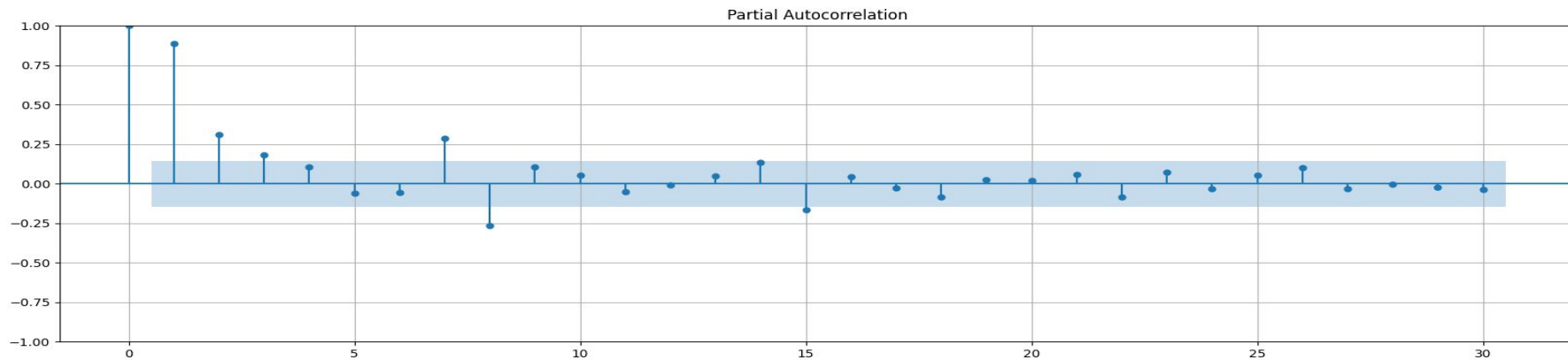
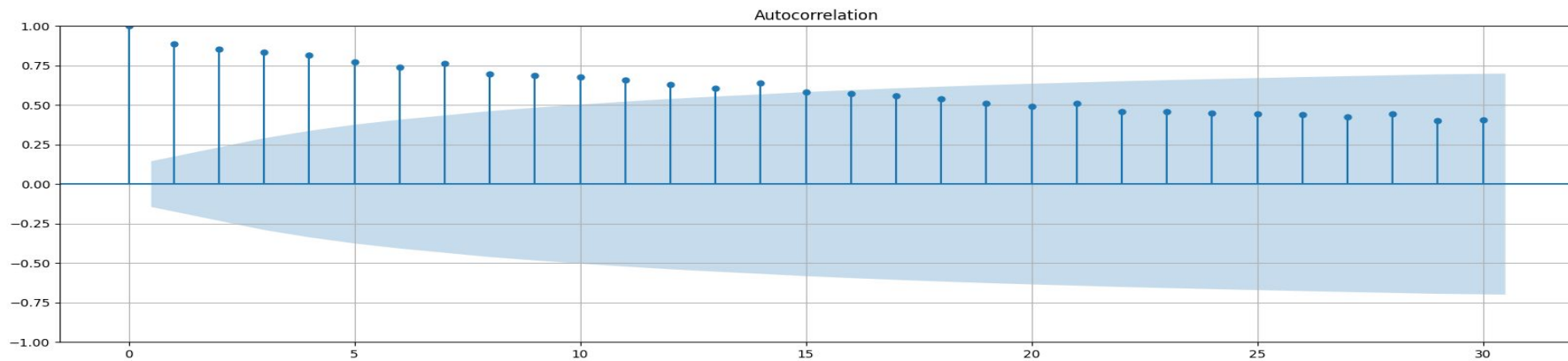
- **Yellow Taxi Trip Records** (PARQUET)
- **Green Taxi Trip Records** (PARQUET)
- **For-Hire Vehicle Trip Records** (PARQUET)
- **High Volume For-Hire Vehicle Trip Records** (PARQUET)

Size(**High Volume For-Hire Vehicle Trip Records**) = (8708453, 14)

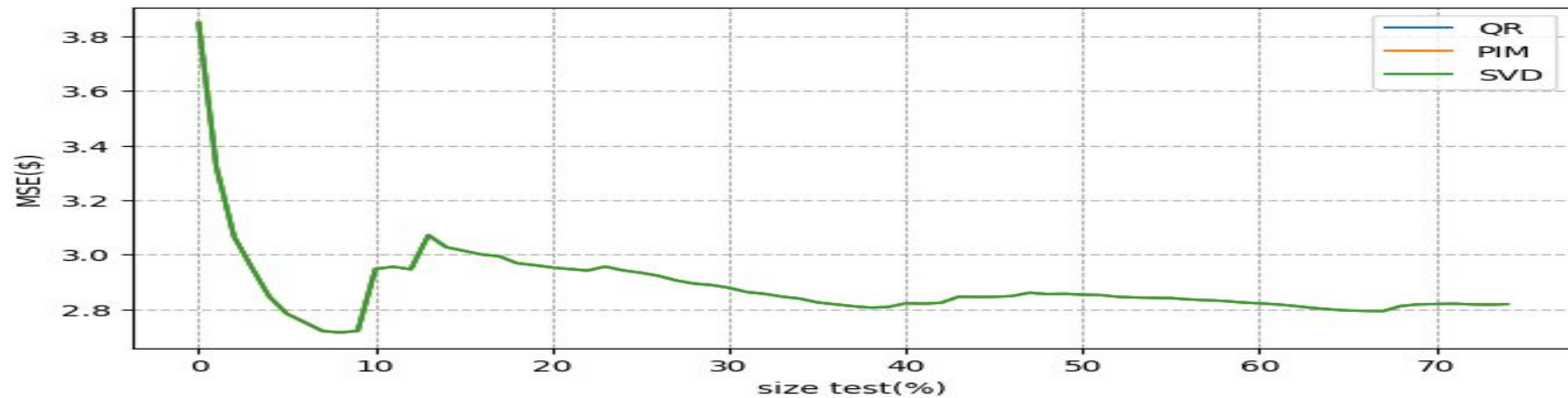
selection of lags



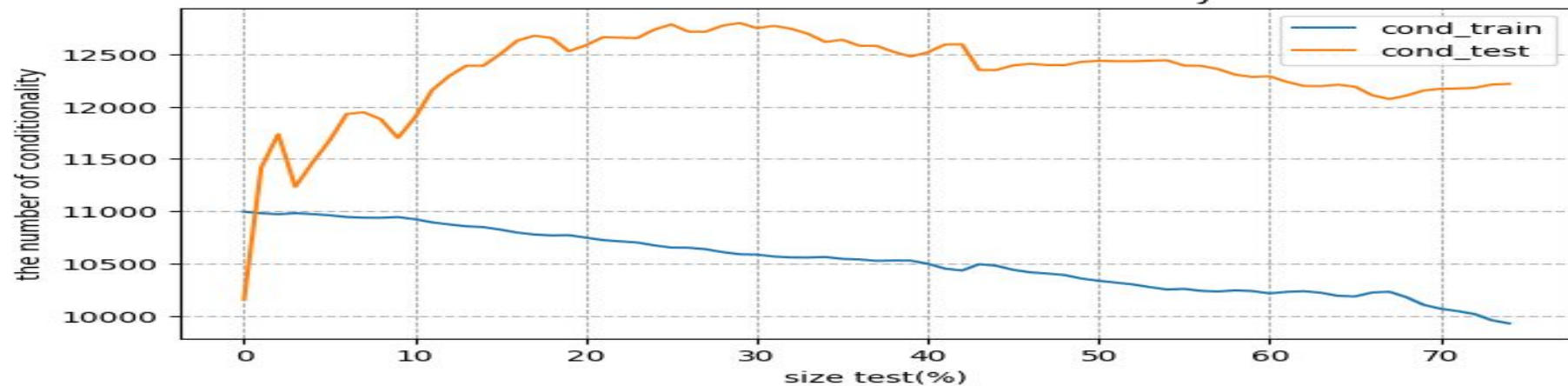
selection of lags



size vs. MSE

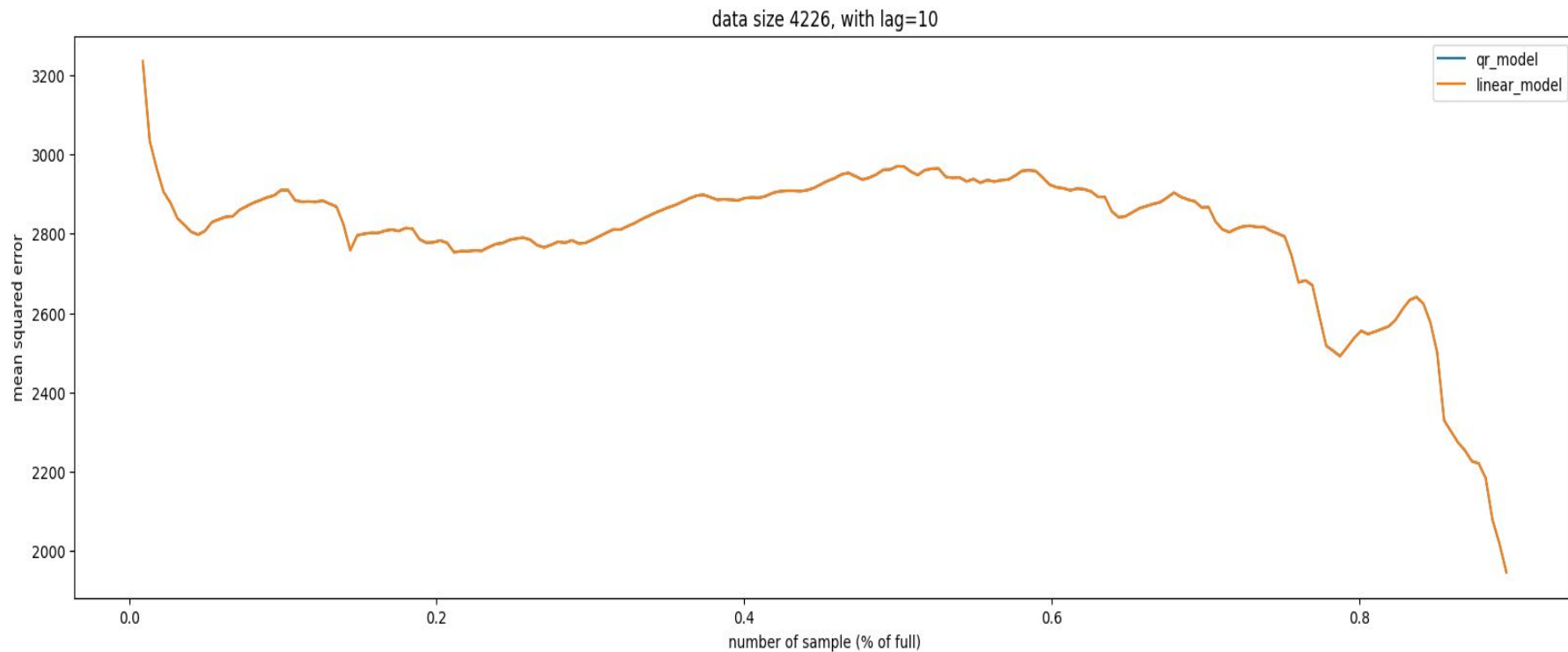


size vs. the number of conditionality

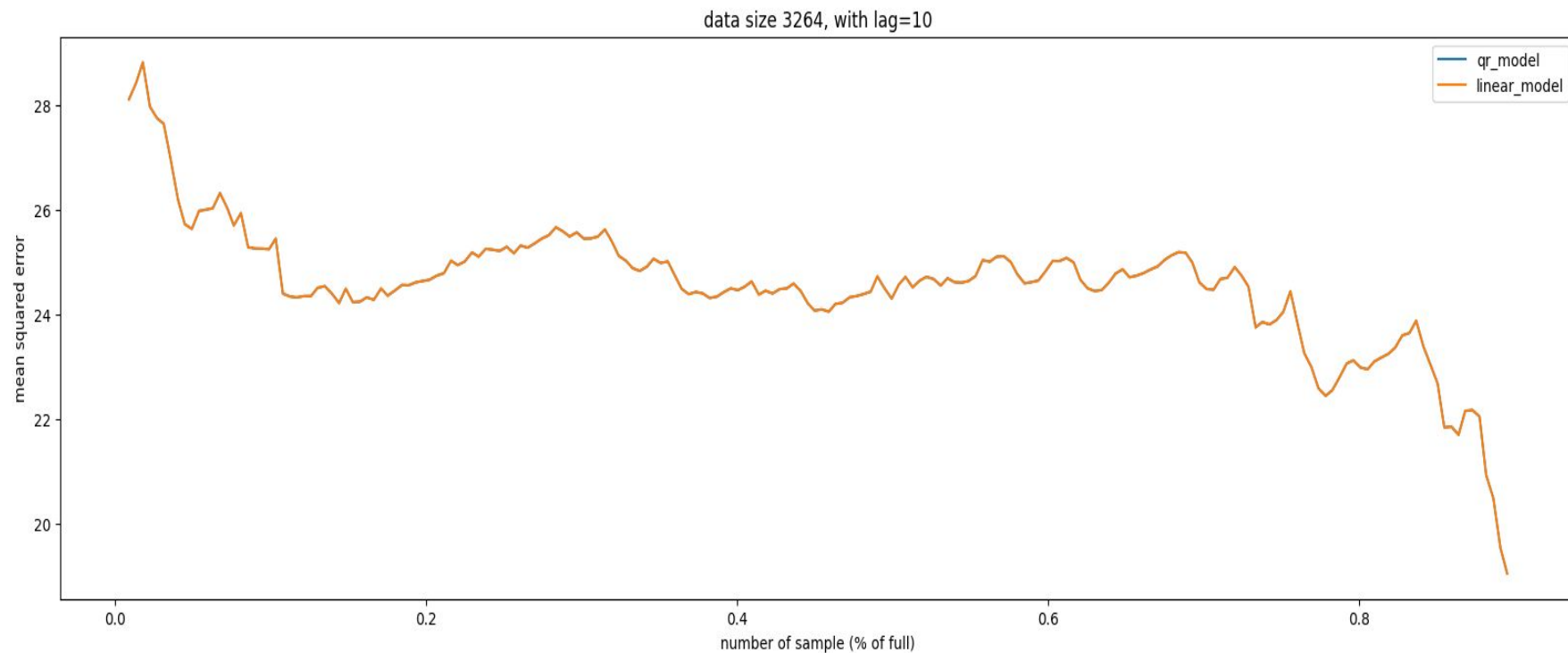


<https://github.com/Mcompetitions/M4-methods/blob/master/Dataset/Train/Daily-train.csv>

V8

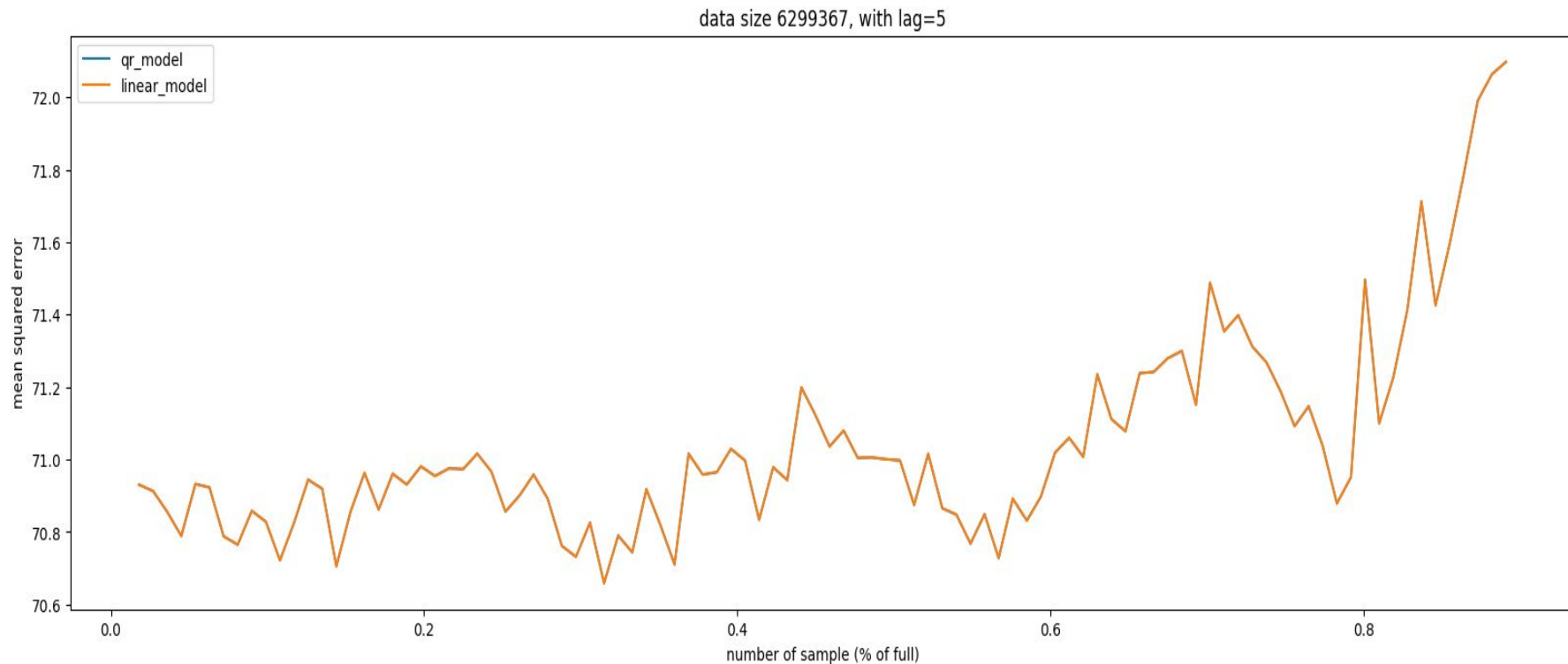


<https://www.kaggle.com/datasets/robervalt/sunspots>



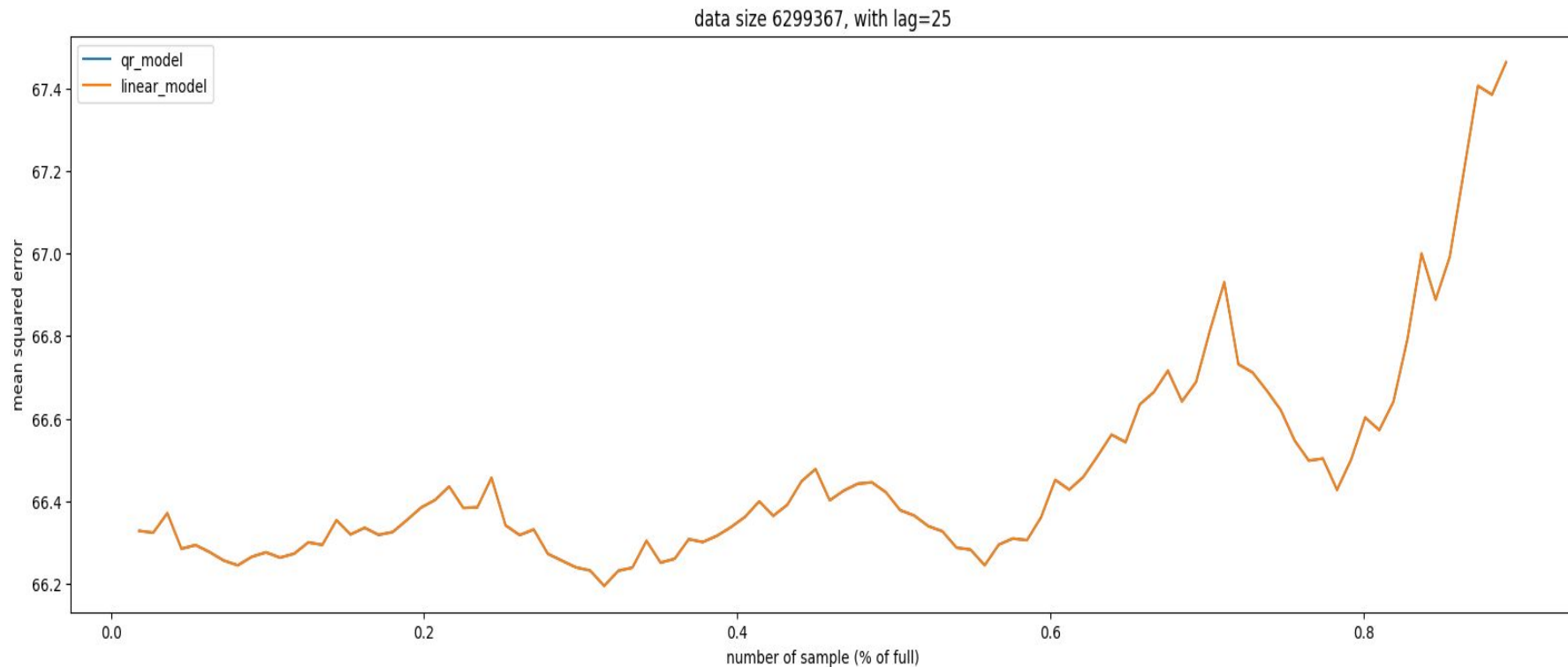
<https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

yellow_tripdata_2020-02.parquet - PULocationID



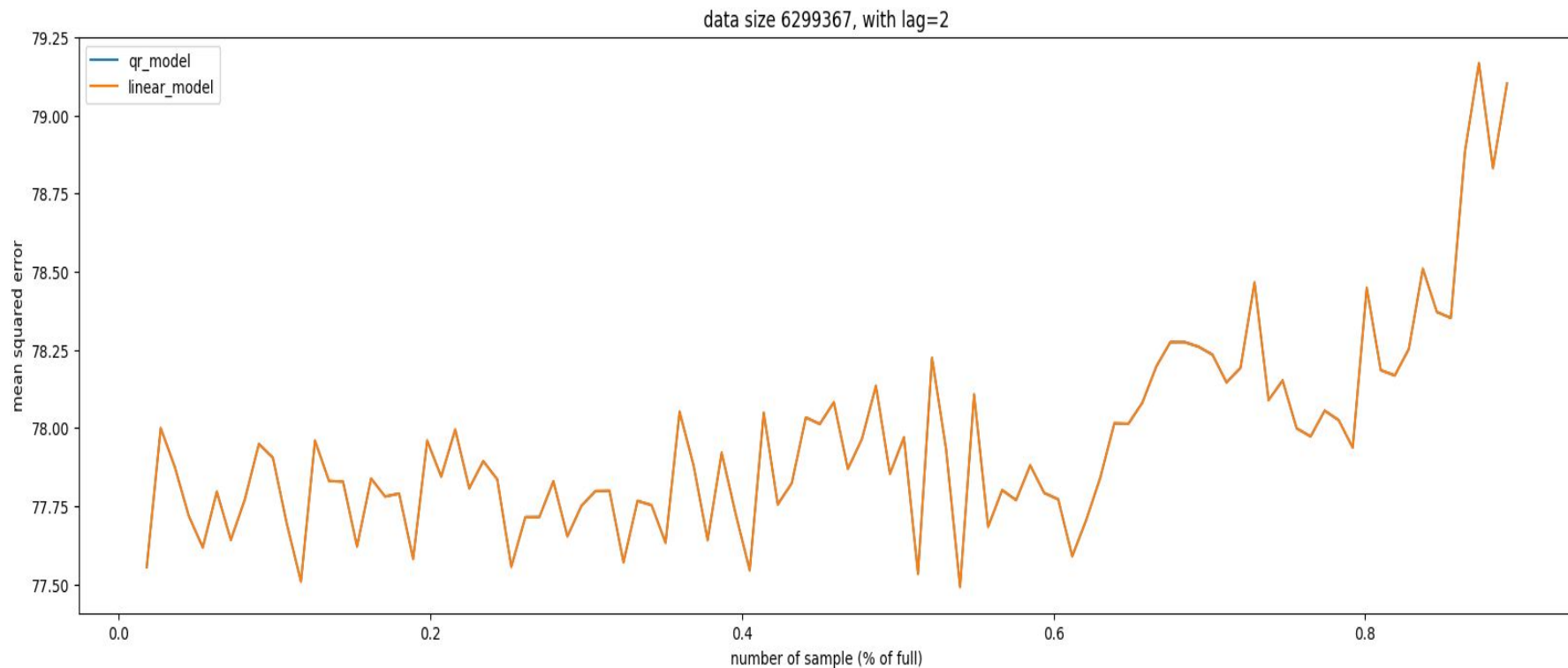
<https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

yellow_tripdata_2020-02.parquet - PULocationID

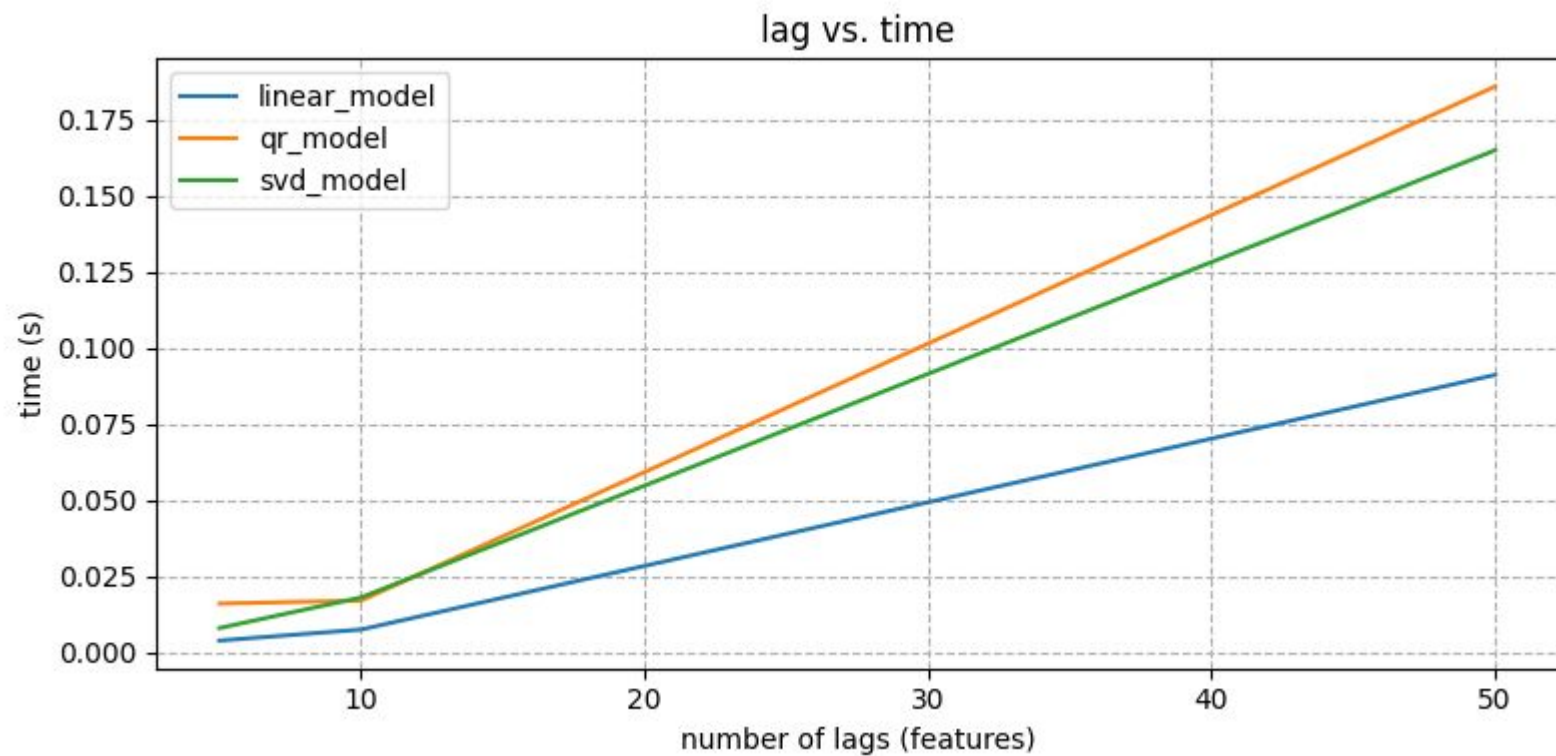


<https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>

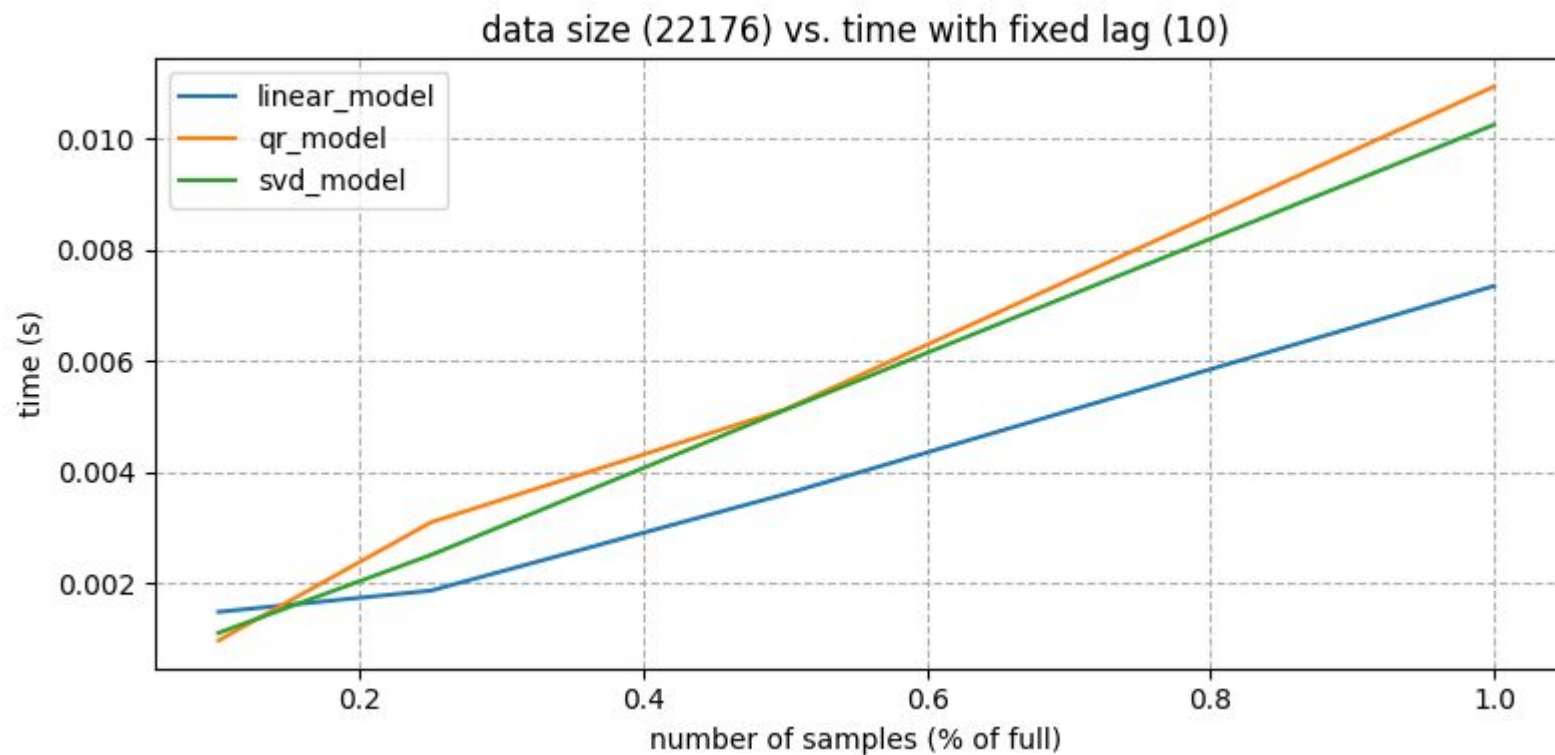
yellow_tripdata_2020-02.parquet - PULocationID



Results



Results



Conclusions

- all the methods are numerically stable (the same solution to SLAE)
- solving via pseudo inverse is faster
- https://github.com/Ulycecc/project_nla

Plans

- Building an accurate taxi payment model
- Find the widest possible class of random processes with a full-rank matrix

https://github.com/Ulycecc/project_nla