# DS-GA 1013 Project Proposal

Andrew Hopen and Andrei Kapustin

For our project, we would like to explore a practical application of the method proposed in the paper "Alternating Minimization for Mixed Linear Regression" by Yi et al [1]. In mixed linear regression, we model some continuous outcome variable as a linear combination of features, with the additional complication that each sample belongs to one of multiple latent classes, and those classes have different regression parameters. Mixed linear models are fit with the EM algorithm: we start by randomly assigning latent membership and randomly initializing the parameter vectors, then we alternate between re-assigning the classes and re-estimating the regression parameters until convergence. Yi et al propose an improvement to the initialization step. Rather than starting with random parameter vectors, they construct a matrix $M$ whose eigenvectors are unbiased estimates of the true parameter vectors (Fig 1a, steps 1-2). Then they search over the unit circle in the space spanned by those eigenvectors for better options, since the true parameter vectors need not be orthogonal (Fig 1a, steps 3-4). Yi et al show that for a constructed dataset, this initialization procedure dramatically speeds up convergence (Fig 1b).



(a) Initialization Algorithm
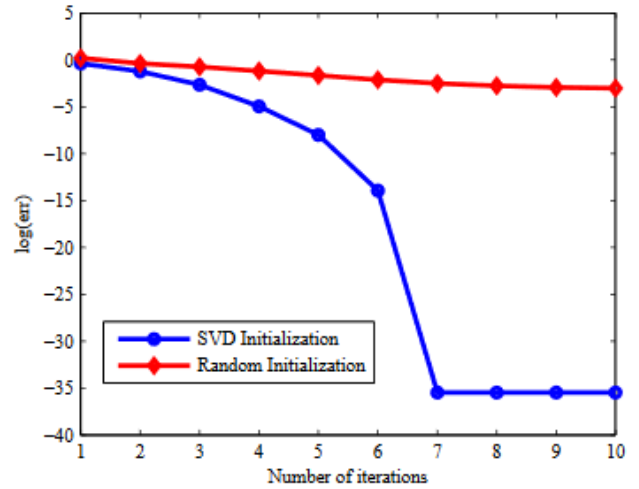


(b) Convergence Speedup

Figure 1: Yi et al Method

We plan to test this procedure on publicly available NYC Taxi Data [2]. This data includes information on all taxi rides in NYC. For example, Figure 2 shows a subset of "green taxi" (outer borough) rides in January 2019. We have color-coded the trips by "RateCode" and shown their corresponding best-fit lines. We see, for example, that trips from the airport JFK have a flat rate of \$52, whereas standard rides show a slope of \$2.29 per mile.

We would like to run the mixed linear regression procedure on this data to see how well we can recover RateCodes and the regression parameters, and to assess speedup when applying the proposed initialization method. We will run these tests for different numbers of covariates and different numbers of latent RateCodes, as permitted by the data.
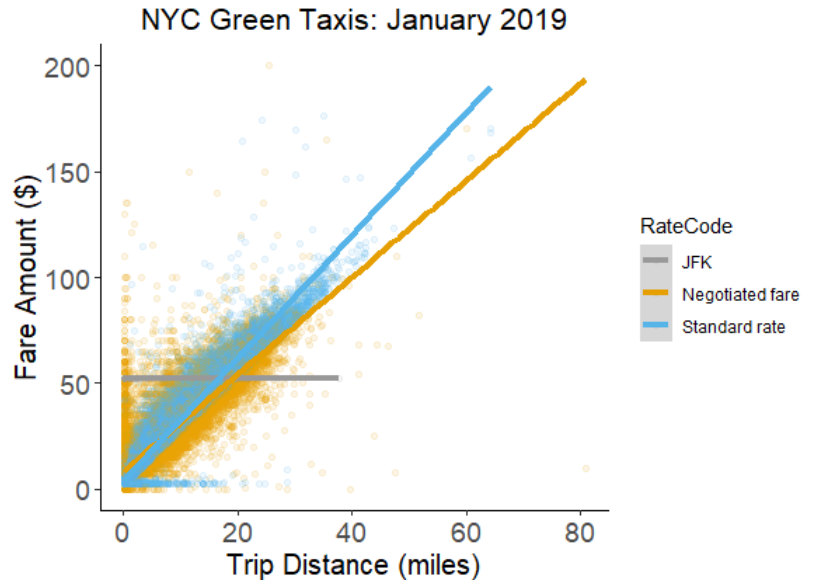


Figure 2: Real-World Mixed Regression

# Bibliography

[1] Xinyang Yi, Constantine Caramanis, Sujay Sanghavi. *Alternating Minimization for Mixed Linear Regression*, 2013
https://arxiv.org/pdf/1310.3745v1.pdf

[2] NYC Taxi Data. https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page