Demographic Model Selection with Deep Learning Ariella L. Gladstein^{1*} and Daniel R. Schrider¹

1) Department of Genetics, University of North Carolina, Chapel Hill, NC aglad@med.unc.edu, https://github.com/agladstein

Abstract: One way to go about obtaining a better understanding of a population's demographic history is to design several competing models and see which best fits the data. We developed a method using a convolutional neural network (CNN) to choose the best demographic model using genome sequence data. We approach demographic model choice as a classification problem, where we train a CNN on simulated genome sequence data from a variety of models and parameter space to learn to recognize the genomic patterns resulting from each type of model. This allows us to input genome sequence from an unknown demographic history, and output posterior probabilities for each model the neural network was trained on. We also designed a feature vector of 30 genomic summary statistics to use in a random forest or approximate Bayesian computation. Overall, the random forest performs the best, while the CNN approaches the accuracy of the random forest with more training data. ABC, on the other hand, was unable to complete all the given classification problems due to linear dependency among the summary statistics.

Motivation

- An important task of demographic inference is model selection.
- Model selection can be used to test evolutionary hypotheses regarding populations' histories.
- Current methods to choose the best model can be grouped into two general categories: likelihood-based and Approximate Bayesian Computation (ABC).
- Likelihood and ABC approaches vary in accuracy, ABC is more flexible than likelihood approaches, and both are computationally burdensome.





- *N_e* parameters ~*U*(1000, 40000)
- Time parameters $\sim U(1, 4000)$ generations
- Migration rate ~U(0.1, 0.9)

