

# Detecting Selection in Experimental Evolution Experiment

Arya Iranmehr  
airanmehr@ucsd.edu

Bafna Lab  
University of California, San Diego

March, 2016

# Introduction

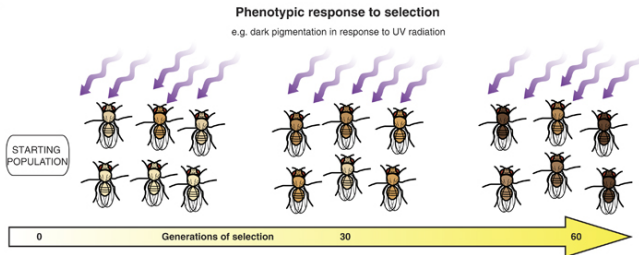
- We are interested in identifying genetic adaptations in different organisms in find genes (or alleles) that are beneficial w.r.t. a **selective pressure**.
- Examples of interesting selective pressures: harsh environmental conditions, antibiotics, chemotherapy, etc.

# History

- Traditionally, detecting/locating selection is done by analysing Allele Frequency Spectrum(AFS) or Haplotypes of a region (say 50Kbp).
- **Now**
  - 1 We can sequence the population at different generations
  - 2 We can't wait until fixation.

# Experimental Evolution

**a**



**b**

**Genotypic response to selection**

e.g. causative variant increases in frequency



**c**

**Pool-Seq**

e.g. base and selected population



# Goals

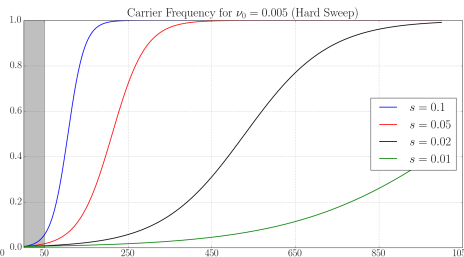
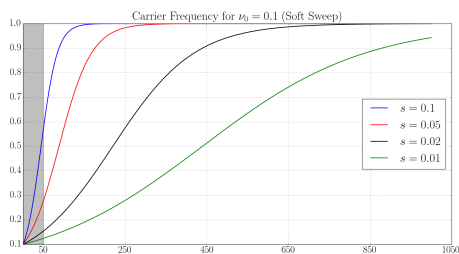
- Given the genome time series data (allele frequencies) we are interested in
  - I Detecting selection
  - II Locating selection
    - (i) Identify gene under selection
    - (ii) Identify the mutation
  - III Estimating model parameters

# Simulations

- (i) For each simulation, population of  $F = 200$  founder lines is created in `msms` program with parameters
  - window size  $L = 50\text{Kbp}$ .
  - mutation rate  $\theta = 200$ .
- (ii) Using  $F = 200$  founder lines a population diploid is created.
- (iii) Using forward simulator population is evolved and AF is sampled every 10 generation for 50 generations.
- (iv) This procedure is repeated 100 times for each  $s = \{0.1, 0.05, 0.02, 0.01, 0\}$ .

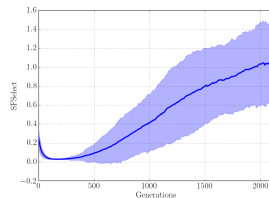
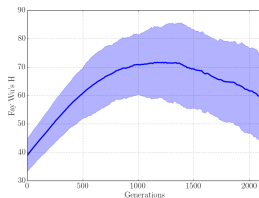
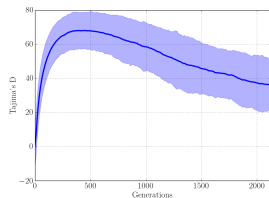
# Single Locus Logistic Model: Carrier Frequency

Dynamics of sweep extremely depend on  $s$  as well as initial carrier frequency. Sampling times are very important prediction performance!



# Bottleneck

- side-effect of not tracking new mutations and restricting population to founder-lines.





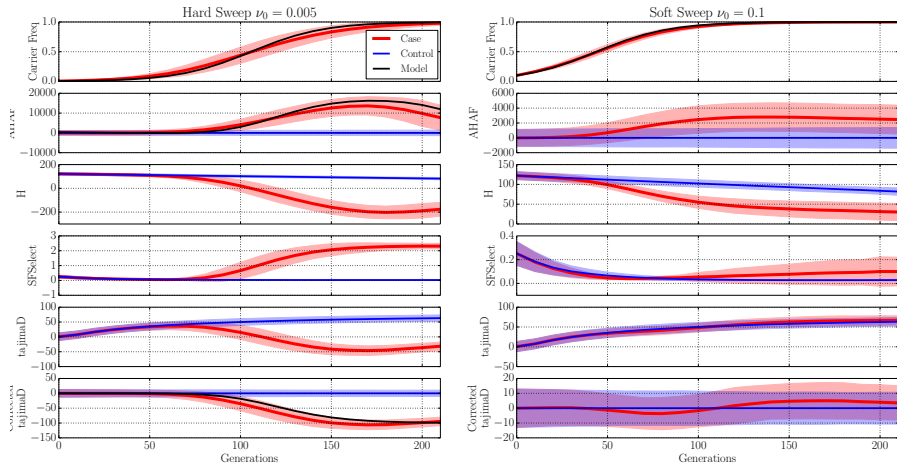


Figure: Mean and 95% CI of 1000 simulations for strong selection.

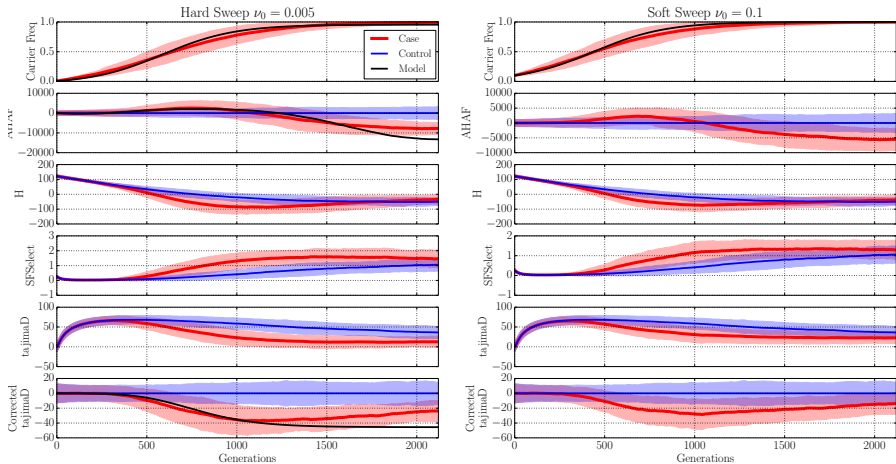
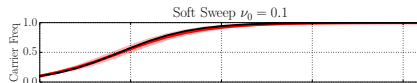
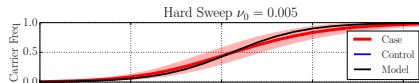


Figure: Mean and 95% CI of 1000 simulations for weak selection.

# Single-Locus Least-Squares Method I

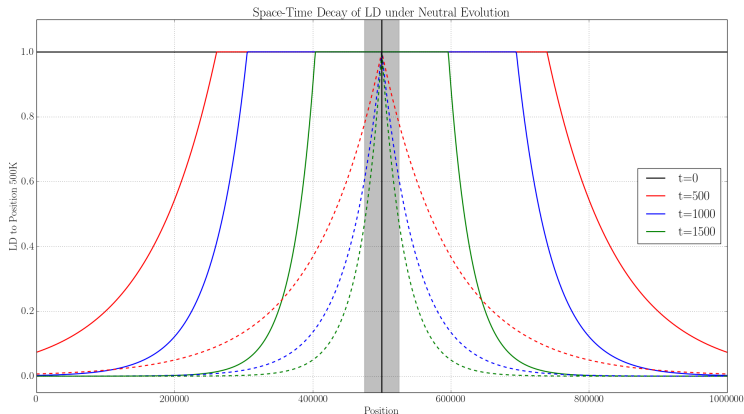
## 1 Model is consistent with observations

Dynamics of Sweep for  $s=0.1$



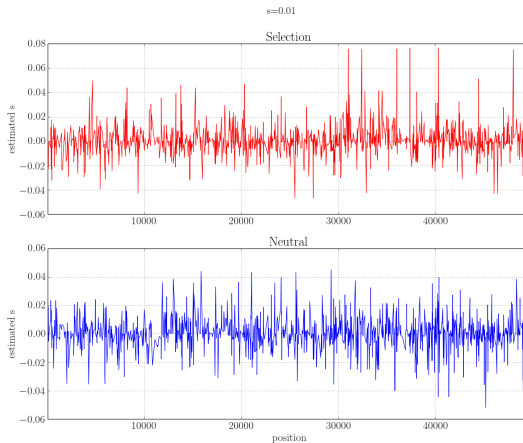
# Single-Locus Least-Squares Method II

## 2 strong linkage in a window



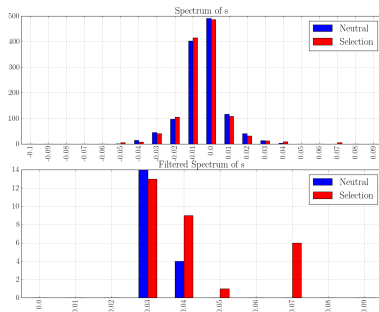
# Single-Locus Least-Squares Method III

## 8 Estimating $s$ at each site



# Single-Locus Least-Squares Method IV

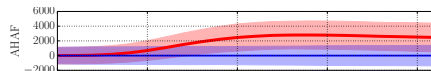
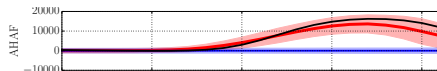
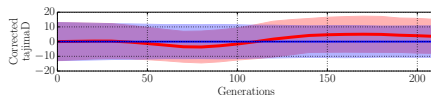
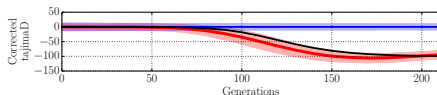
- 4 spectrum of  $\hat{s}$  for a window



- 5 We considered avg of top  $\hat{s}$  as a predictor of selection for a window.

# Tajima's D and Fay Wu's H

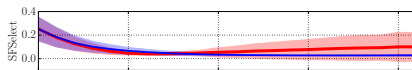
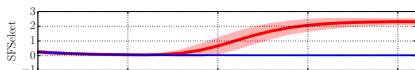
- 1 After correcting observations for bottleneck effect the model look like observations, but it has a higher variance.



- 2 The parametric models assumes hard-sweep, so the model become invalid in soft sweep.

# SFSelect and Gaussian Process

- 1 SFSelect process is a monotone process, so we summed all the values in the time series to get a predictor for selection.



- 2 GP is the state-of-the-art model, proposed by Terhorst et al, which fits a Gaussian process on the time series data.



# Experiments

- for each setting we performs 200 simulations, (100 neutral and 100 selection), and predictive performance of all the methods in detecting selection.
- We computed ROC curve and defined power of a method as area under ROC curve when False-Positive rate is less than 0.1.



# Summary

- Behaviour of different population statistics in time is studied.
- A method based on single locus AF is proposed and shown to have superior performance than multi-locus and traditional methods.
- Analysis of real data and locating genes in the genome is our next step.

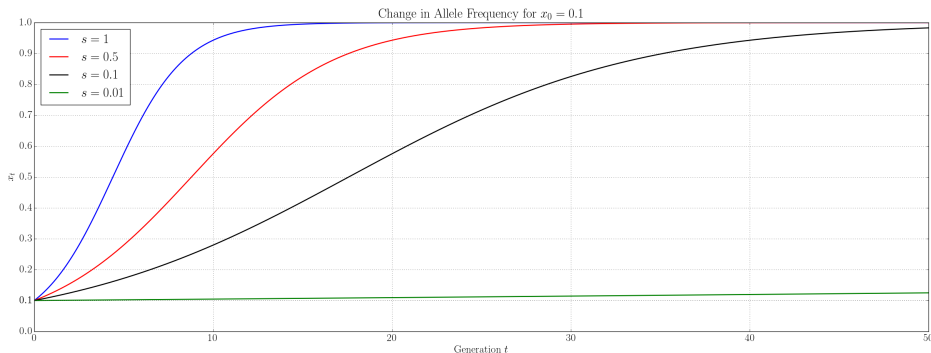
Thanks!

# Single Locus Model: Carrier Frequency

- By differentiating update equations ( $x_{t+1} = x_t + \frac{sx_t(1-x_t)}{2+2sx_t}$  w.r.t.  $t$  and solving differential equation, we have

$$\nu_t = \sigma(st/2 + \eta(\nu_0))$$

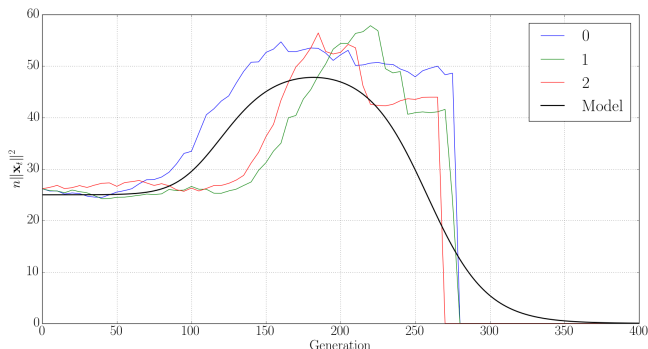
where  $\sigma(\cdot)$  is logistic function and  $\eta = \sigma^{-1}$  is the logit function and  $\nu_t$  is the frequency of the carrier at time  $t$ .



# Multi Locus Model: Average Haplotype Allele Frequency

$$\mathbb{E}[1\text{-HAF}(t)] = \|\mathbf{x}_t\|^2 \approx \theta \nu_t \left( \frac{\nu_t + 1}{2} - \frac{1}{(1 - \nu_t)n + 1} \right) + \theta(1 - \nu_t) \left( \frac{n + 1}{2n} - \frac{1}{(1 - \nu_t)n + 1} \right)$$

where  $\mathbf{x}_t$  is vector of AF at time  $t$  and  $\nu_t = \sigma(st/2 + \eta(\nu_0))$ .



- Its too complex, likely to overfit with small number of iid replicates.
- Although the likelihood model is based on different parameters, in practice, it can learn only one parameter at a time.
- not tractable, its time complexity is quartic!
- worse, each iteration requires **maxGeneration** recursion which makes it very hard to analyse late epochs of sweeps.
- In addition to PoolSeq data, it requires initial population haplotypes.
- Despite its elegant theory, it has not compared with classical methods.
- In practice, single locus scan is performed and multi-locus (with 3-7 seg. sites.) model is fitted at regions of interests.