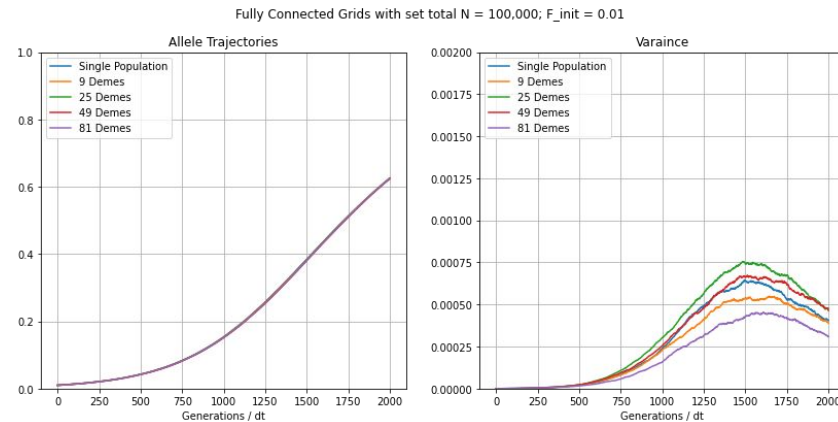


Rotation Week 8

Updates



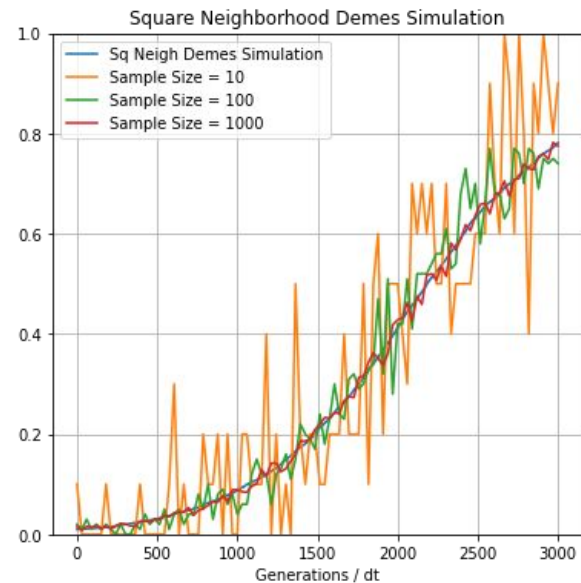
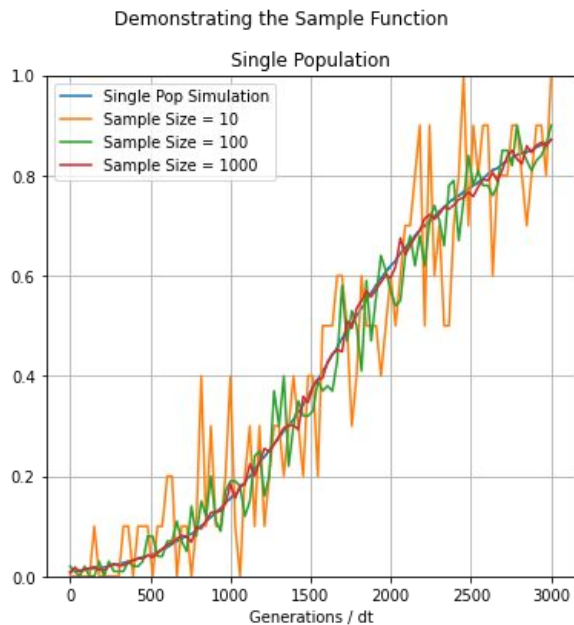
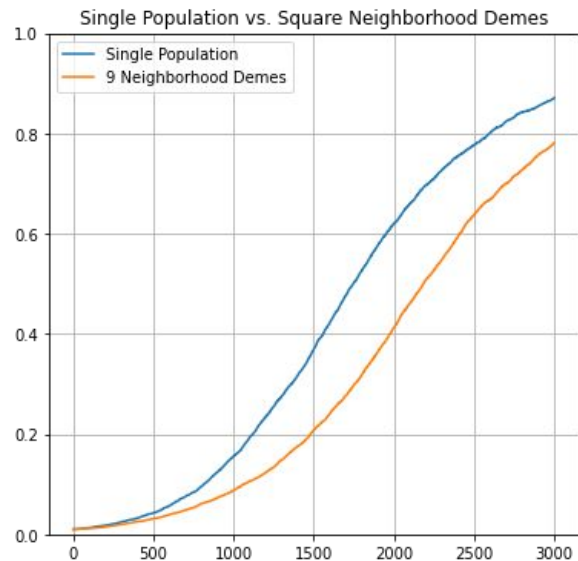
Rowan Hart
Ecology and Evolution
University of Chicago



Sampling Data from Simulations

```
##### Random Sampling From Simulations #####  
  
# Function to Sample From Simulation  
def sample_sim(sim, num_samples, samp_size):  
    sample_times = np.linspace(0, len(sim) - 1, num_samples, dtype=int)  
    freqs = np.zeros((num_samples, samp_size))  
    inds = np.zeros((num_samples, samp_size))  
    samps = np.zeros(num_samples)  
    for i in np.arange(num_samples):  
        for j in np.arange(samp_size):  
            freqs[i,j] = np.random.choice(sim[sample_times[i]])  
            inds[i,j] = np.random.binomial(1, freqs[i,j])  
        samps[i] = sum(inds[i]) / samp_size  
    return (sample_times, samps)
```

Sampling Data from Simulations



What are the probabilities of **allele loss**?

Differences between Normal and Binomial Distributions

Binomial Distribution:

$$\text{CDF}(\mathbf{X} = 0, 2*N, \text{freq})$$

Normal Distribution:

$$\text{CDF}(\mathbf{X} = 0, \text{mean}=\text{freq}, \text{std}=(1/2N)*(1-\text{freq})*\text{freq})$$

Parameters:

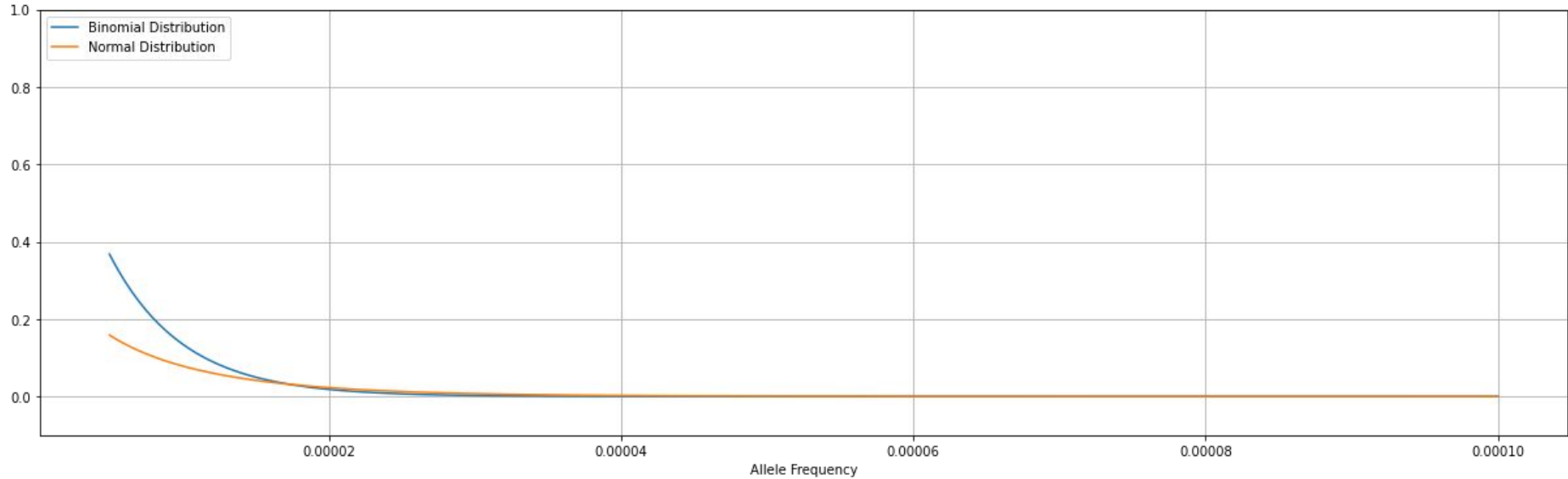
N_total = 100,000; No Selection

- Test along a number of allele frequencies

What are the probabilities of **allele loss**?

Differences between Normal and Binomial Distributions

The probability of an allele getting lost after 1 Generation



What are the probabilities of **allele loss**?

Differences between Normal and Binomial Distributions

Binomial Distribution:

$$N_i = N / d$$
$$[\text{CDF}(\mathbf{X} = 0, 2*N_i, \text{freq})]^d$$

Normal Distribution:

$$N_i = N / d$$
$$[\text{CDF}(\mathbf{X} = 0, \text{mean}=\text{freq}, \text{std}=(1/2N_i)*(1-\text{freq})*\text{freq})]^d$$

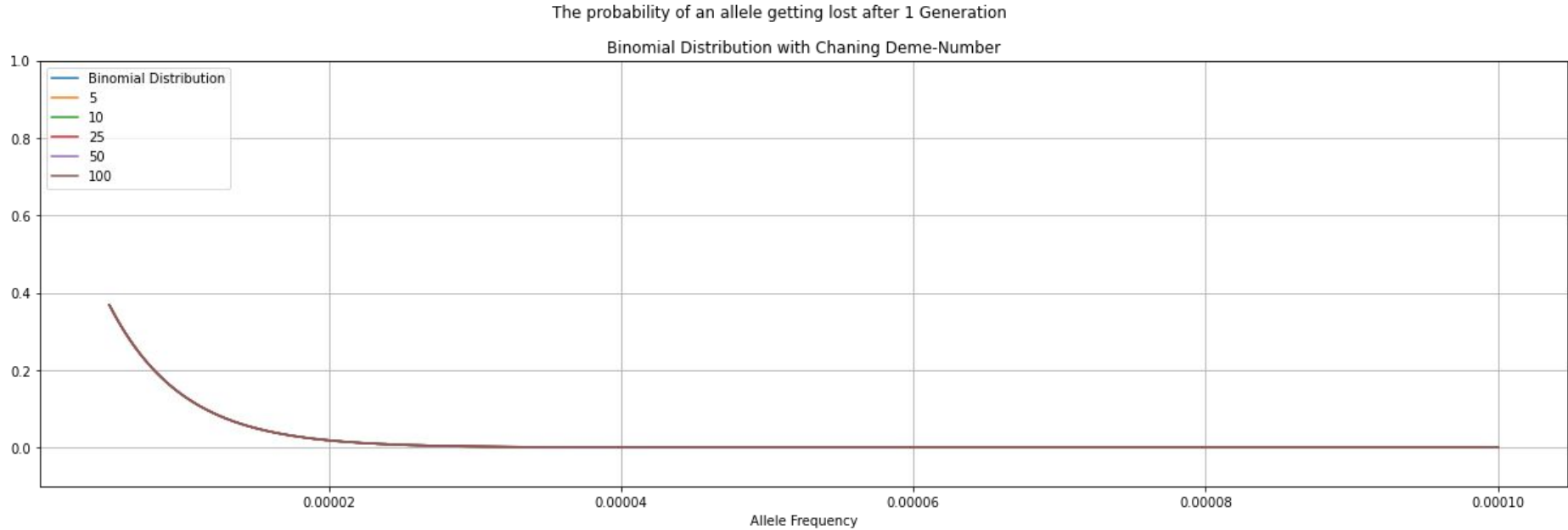
Parameters:

$N_{\text{total}} = 100,000$; No Selection

- Test along a number of allele frequencies
- Test along $d = [1, 5, 10, 25, 50, 100]$

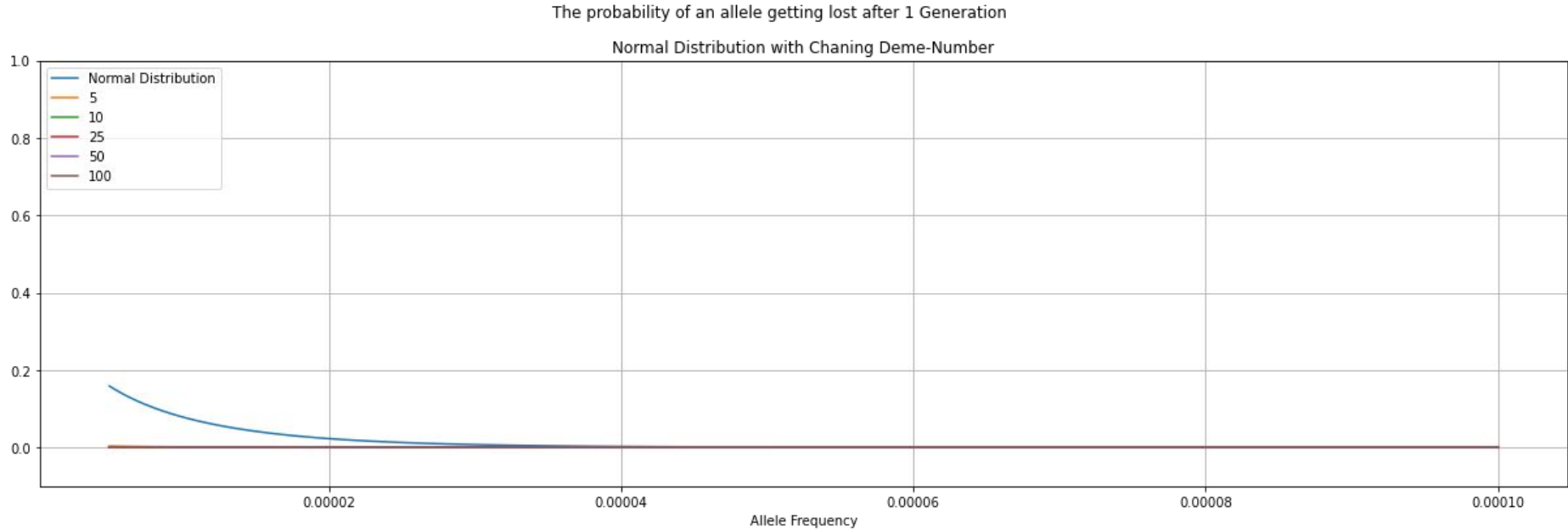
What are the probabilities of **allele loss**?

Differences between Normal and Binomial Distributions



What are the probabilities of **allele loss**?

Differences between Normal and Binomial Distributions



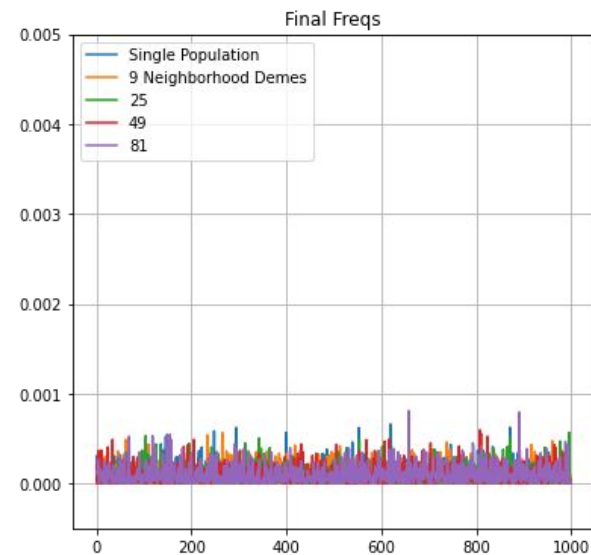
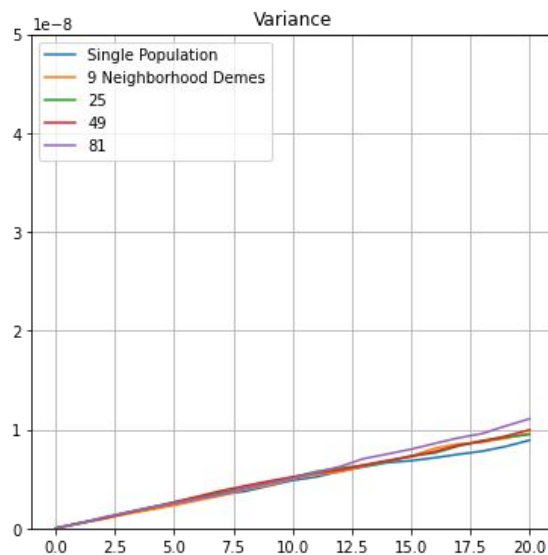
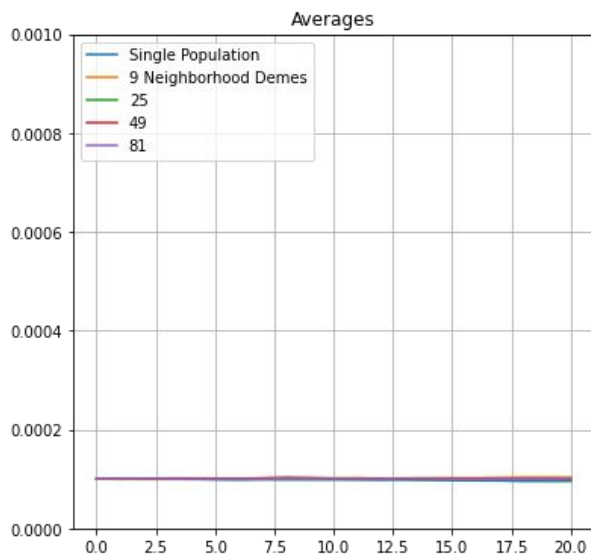
Comparing **Binomial Distribution** to Gaussian distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.00005$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Binomial Simulations: Averages w/ no selection and $F_{\text{init}} = .0001$



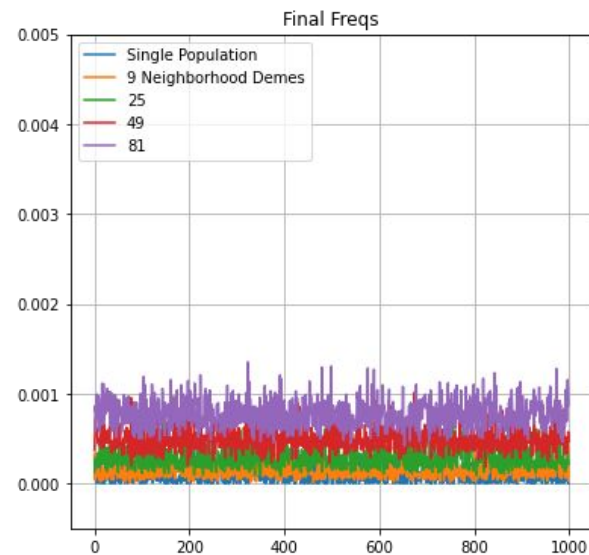
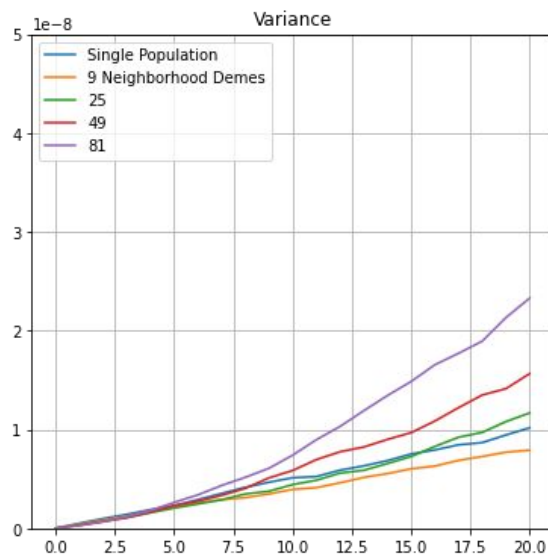
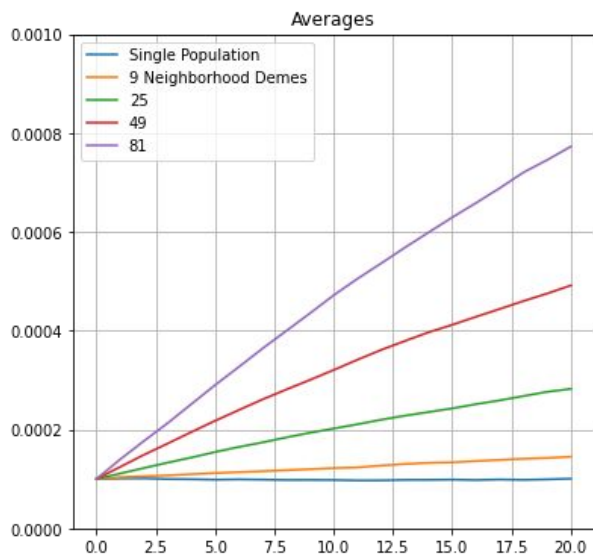
Comparing Binomial Distribution to **Gaussian distribution**

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.0001$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Normal Simulations: Averages w/ no selection and $F_{\text{init}} = .0001$



Comparing **Gaussian distribution** to Binomial Distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.00005$; $dt = 1$; $t = 20$ generations

- no selection
- fully connected demes

Binomial Alleles Lost
(of 1,000 iterations to 20 gen)

```
print(sum(results_bin[0][-1] < (10 ** -15)))  
print(sum(results_bin[1][-1] < (10 ** -15)))  
print(sum(results_bin[2][-1] < (10 ** -15)))  
print(sum(results_bin[3][-1] < (10 ** -15)))  
print(sum(results_bin[4][-1] < (10 ** -15)))
```

634
615
595
643
638

Normal Alleles Lost
(of 1,000 iterations to 20 gen)

```
print(sum(results[0][-1] < (10 ** -15)))  
print(sum(results[1][-1] < (10 ** -15)))  
print(sum(results[2][-1] < (10 ** -15)))  
print(sum(results[3][-1] < (10 ** -15)))  
print(sum(results[4][-1] < (10 ** -15)))
```

349
0
0
0
0

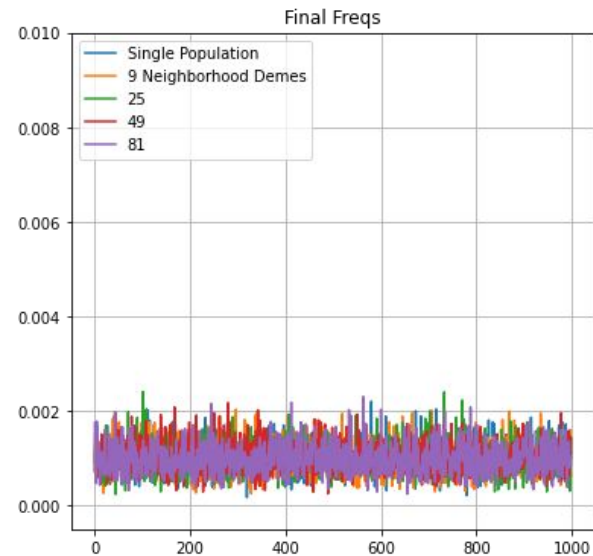
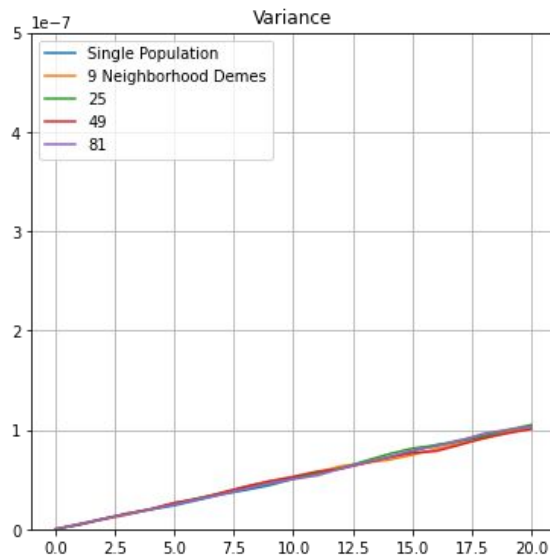
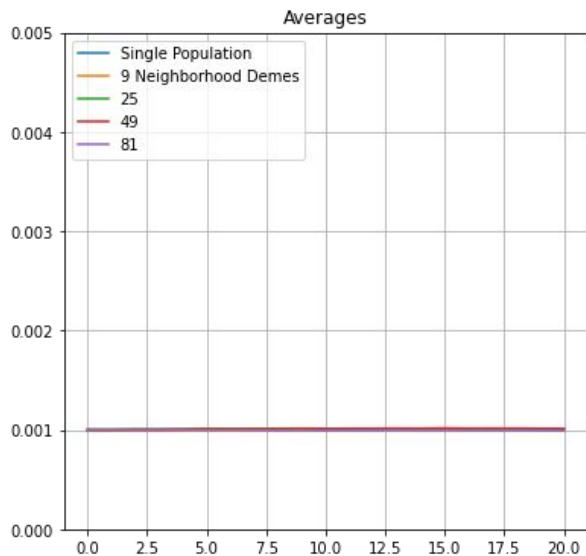
Comparing **Binomial Distribution** to Gaussian distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.001$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Binomial Simulations: Averages w/ no selection and $F_{\text{init}} = .001$



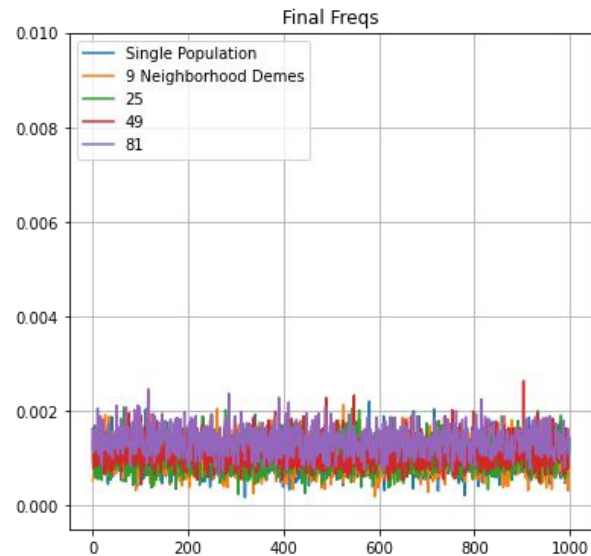
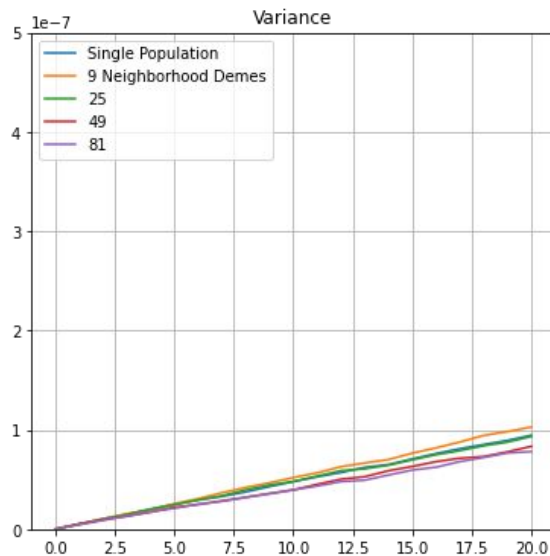
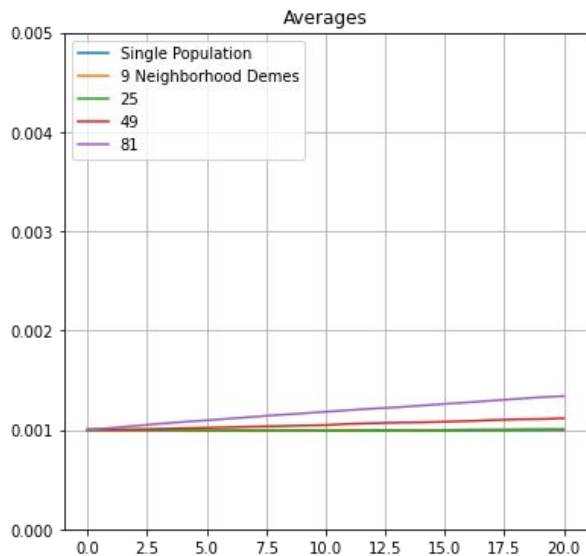
Comparing Binomial Distribution to **Gaussian** distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.001$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Normal Simulations: Averages w/ no selection and $F_{\text{init}} = .001$



Comparing **Gaussian distribution** to Binomial Distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.001$; $dt = 1$; $t = 20$ generations

- no selection
- fully connected demes

Binomial Alleles Lost
(of 1,000 iterations to 20 gen)

```
print(sum(results_bin[0][-1] < (10 ** -15)))  
print(sum(results_bin[1][-1] < (10 ** -15)))  
print(sum(results_bin[2][-1] < (10 ** -15)))  
print(sum(results_bin[3][-1] < (10 ** -15)))  
print(sum(results_bin[4][-1] < (10 ** -15)))
```

0
0
0
0
0

Normal Alleles Lost
(of 1,000 iterations to 20 gen)

```
print(sum(results[0][-1] < (10 ** -15)))  
print(sum(results[1][-1] < (10 ** -15)))  
print(sum(results[2][-1] < (10 ** -15)))  
print(sum(results[3][-1] < (10 ** -15)))  
print(sum(results[4][-1] < (10 ** -15)))
```

0
0
0
0
0

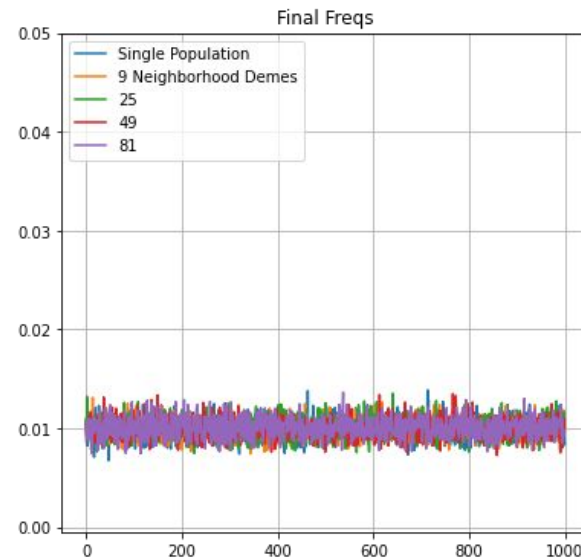
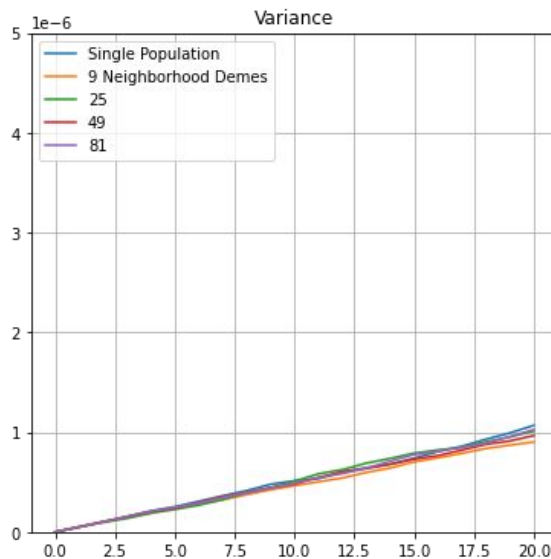
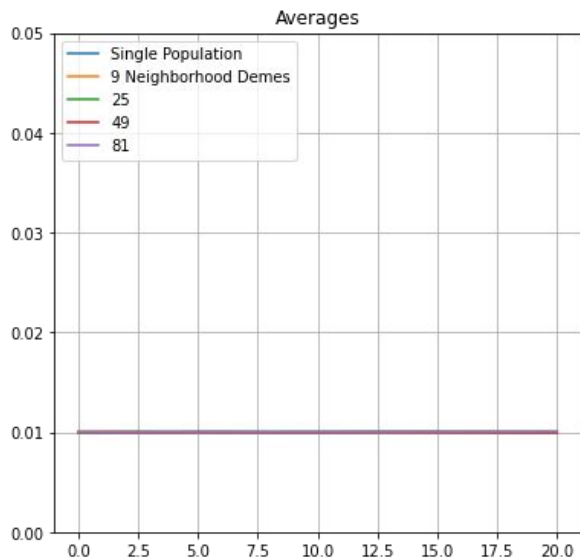
Comparing **Binomial Distribution** to Gaussian distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.01$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Binomial Simulations: Averages w/ no selection and $F_{\text{init}} = .001$



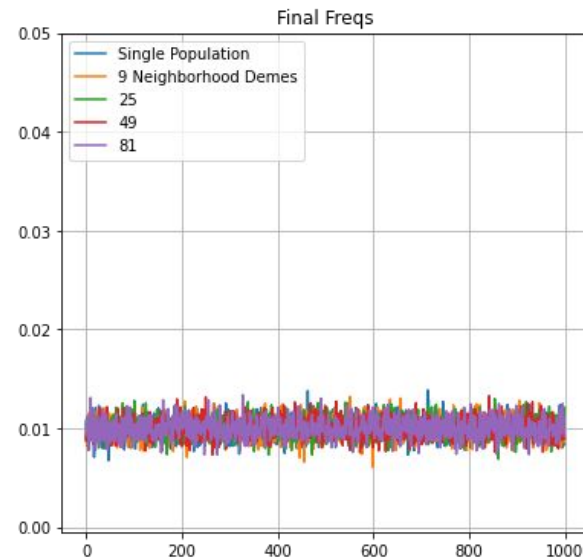
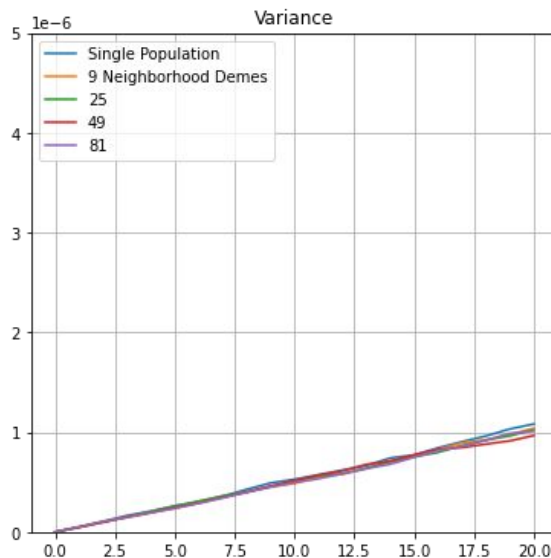
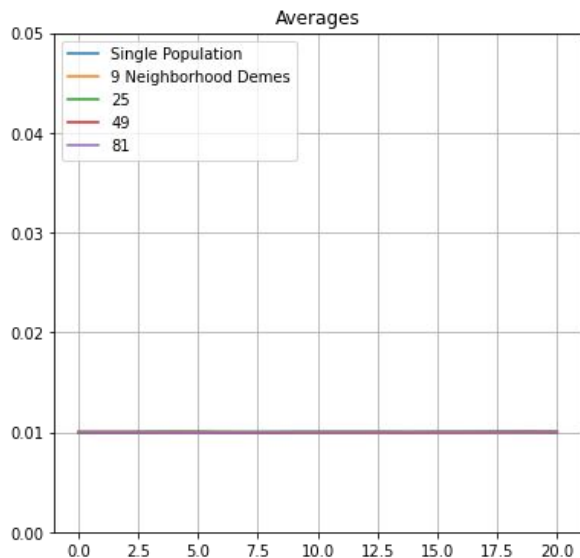
Comparing Binomial Distribution to **Gaussian** distribution

Parameters:

$N_{\text{total}} = 100,000$; $F_{\text{init}} = 0.01$; $dt = 1$; $t = 20$ generations; 1,000 iterations

- no selection
- fully connected demes

Testing Normal Simulations: Averages w/ no selection and $F_{\text{init}} = .001$



Inference of Population Structure from Time-Series Genotype Data

Tyler A. Joseph¹  , Itzik Pe'er^{1 2 3}  

We assume each sampled individual is a mixture of K unobserved populations. Let the vector $\theta_d = (\theta_{d1}, \theta_{d2}, \dots, \theta_{dK})$ give the proportion of the genome individual d inherits from each population k ; hence $\sum_{k=1}^K \theta_{dk} = 1$. Denote the allele frequency of the reference allele at locus l in population k at time point t by β_{kl}^t . With this notation, we assume the following generative model:

$$\theta_d \sim \text{Dirichlet}(\alpha_1, \alpha_2, \dots, \alpha_K) \quad \text{for } d = 1, \dots, D. \quad (\text{Equation 1})$$

$$\beta_{kl}^t | \beta_{kl}^{t-1} \sim \text{Normal}\left(\beta_{kl}^{t-1}, \frac{g_t - g_{t-1}}{12N_k}\right) \quad \text{for } l = 1 \dots L \text{ and } K = 1 \dots K \quad (\text{Equation 2})$$

$$x_{dl} | \beta_{1:K,l}^t, \theta_d \sim \text{Binomial}\left(n_d, \sum_{k=1}^K \theta_{dk} \beta_{kl}^t\right) \quad \text{for } d = 1 \dots D \text{ and } l = 1 \dots L \quad (\text{Equation 3})$$

with initial allele frequencies β_{kl}^0 and effective population sizes N_k treated as parameters. The β_{kl}^t are estimated from data, while N_k are treated as known and fixed. Ancient DNA samples are typically pseudo-haploid, as low sequencing depth makes it difficult to call full diploid genotypes. We explicitly model pseudo-haploid individuals by setting the sample size parameter of the binomial as either $n_d = 2$ or $n_d = 1$ depending on whether an individual is diploid or pseudo-haploid.

Note that the variance for drift is different than the variance obtained under Wright-Fisher or using a diffusion approximation. Traditionally, the variance of $\beta_{kl}^t | \beta_{kl}^{t-1}$ is

$$\text{Var}(\beta_{kl}^t | \beta_{kl}^{t-1}) = \frac{\beta_{kl}^{t-1}(1 - \beta_{kl}^{t-1})(g_t - g_{t-1})}{2N_k}.$$

However, the appearance of allele frequencies in the variance leads to difficulties in deriving an inference algorithm. We approximate the variance by taking the average variance over possibly allele frequencies:

$$\text{Var}(\beta_{kl}^t | \beta_{kl}^{t-1}) \approx \frac{g_t - g_{t-1}}{12N_k} = \int_0^1 \frac{\beta_{kl}^{t-1}(1 - \beta_{kl}^{t-1})(g_t - g_{t-1})}{2N_k} d\beta_{kl}^{t-1}.$$

Intuitively, this is similar to assuming a uniform prior over β_{kl}^{t-1} in the variance and taking the expectation.