

(Question1): (15 points) The Hardy-Weinberg formula for the genetic variation of a population at equilibrium states that with probability π of the A allele, the probabilities of genotypes (AA, Aa, aa) are $(\pi^2, 2\pi(1-\pi), (1-\pi)^2)$. Assume a multinomial distribution for n observations with counts (y_1, y_2, y_3) of (AA, Aa, aa).

- Express the log-likelihood function for the multinomial and find the MLE of π , denoted as $\hat{\pi}$.
- Find $I(\pi)$, and thus specify the asymptotic distribution of $\hat{\pi}$ (Hint: you may need $\text{Cov}[Y_i, Y_j] = -np_i p_j$ if $(Y_1, Y_2, \dots, Y_K) \sim \text{multinomial}(n, p_1, p_2, \dots, p_K)$).
- Explain how you could use the asymptotic distribution in (b) to construct a 95% confidence interval for π .

(a) The maximum likelihood function is:

$$\begin{aligned} L(\pi) &= (\pi^2)^{y_1} \cdot [2\pi(1-\pi)]^{y_2} \cdot [(1-\pi)^2]^{y_3} \\ &= \pi^{2y_1} \cdot (2\pi)^{y_2} (1-\pi)^{y_2} (1-\pi)^{2y_3} \\ &= 2^{y_2} \cdot \pi^{2y_1+y_2} (1-\pi)^{y_2+2y_3} \end{aligned}$$

$$\ell(\pi) = \log L = y_2 \log 2 + (2y_1 + y_2) \log \pi + (y_2 + 2y_3) \log(1-\pi)$$

$$\begin{aligned} \ell'(\pi) &= \frac{2y_1 + y_2}{\pi} - \frac{y_2 + 2y_3}{1-\pi} \\ &= \frac{2y_1 + y_2 - (y_2 + 2y_3)\pi}{\pi(1-\pi)} \\ &= \frac{2y_1 + y_2 - 2n\pi}{\pi(1-\pi)} \end{aligned}$$

$$\text{Let } \ell'(\hat{\pi}) = 0 \Rightarrow \hat{\pi} = \frac{2y_1 + y_2}{2n}$$

$$\ell''(\pi) = \frac{-2y_1 - y_2}{\pi^2} - \frac{y_2 + 2y_3}{(1-\pi)^2} < 0, \forall \pi$$

$$\text{So } \hat{\pi} = \frac{2y_1 + y_2}{2n} \text{ is the global maximizer of } L(\pi).$$

$$(b) E(y_1) = n \cdot \pi^2$$

$$E(y_2) = 2n\pi(1-\pi)$$

y_1 follows Binomial distribution: $y_1 \sim B(n, \pi^2)$,

$$\text{So } \text{Var}(y_1) = n\pi^2(1-\pi^2).$$

y_2 follows Binomial distribution: $y_2 \sim B(n, 2\pi(1-\pi))$,

$$\text{So } \text{Var}(y_2) = 2n\pi(1-\pi)(1-2\pi+2\pi^2)$$

$$\Rightarrow E(2y_1 + y_2) = 2n\pi^2 + 2n\pi(1-\pi) = 2n\pi$$

$$\begin{aligned} E(2y_1 + y_2)^2 &= \text{Var}(2y_1 + y_2) + [E(2y_1 + y_2)]^2 \\ &= 4\text{Var}(y_1) + \text{Var}(y_2) + 4\text{Cov}(y_1, y_2) + 4n^2\pi^2 \\ &= 4\text{Var}(y_1) + \text{Var}(y_2) + 2\text{Cov}(2y_1, y_2) + 4n^2\pi^2 \\ &= 4n\pi^2(1-\pi^2) + 2n\pi(1-\pi)(1-2\pi+2\pi^2) - 8n\pi^3(1-\pi) + 4n^2\pi^2 \\ &= 2n\pi - 2n\pi^3 + 4n^2\pi^2 \end{aligned}$$

$$\begin{aligned}
 I(\pi) &= E\left(\frac{\partial \ell(\pi)}{\partial \pi}\right)^2 = E\left(\frac{zy_1 + y_2 - 2n\pi}{\pi(1-\pi)}\right)^2 \\
 &= \frac{1}{\pi^2(1-\pi)^2} [E(zy_1 + y_2)^2 - 4n\pi E(zy_1 + y_2) + 4n^2\pi^2] \\
 &= \frac{1}{\pi^2(1-\pi)^2} (2n\pi - 2n\pi^2) \\
 &= \frac{2n}{\pi(1-\pi)}
 \end{aligned}$$

$$\Rightarrow \frac{1}{I(\pi)} = \frac{\pi(1-\pi)}{2n}$$

The asymptotic distribution of $\hat{\pi}$ is $N(\pi, \frac{1}{I(\pi)})$

$$\Rightarrow \hat{\pi} \sim N\left(\pi, \frac{\pi(1-\pi)}{2n}\right).$$

$$(c) \quad 1-\alpha = 95\% \Rightarrow \alpha = 5\% = 0.05 \Rightarrow \frac{\alpha}{2} = 0.025.$$

$$\text{Since } \hat{\pi} \sim N\left(\pi, \frac{\pi(1-\pi)}{2n}\right),$$

$$\frac{\hat{\pi} - \pi}{\sigma} \sim N(0,1), \text{ where } \sigma = \sqrt{\frac{\pi(1-\pi)}{2n}}$$

$$P(-z_{0.025} < \frac{\hat{\pi} - \pi}{\sigma} < z_{0.025}) = 95\%$$

$$\Rightarrow \hat{\pi} - \sigma z_{0.025} < \pi < \hat{\pi} + \sigma z_{0.025}, \text{ where } z_{0.025} = 1.96$$

$$\pi \in \left[\frac{zy_1 + y_2}{2n} - 1.96 \frac{\pi(1-\pi)}{2n}, \frac{zy_1 + y_2}{2n} + 1.96 \frac{\pi(1-\pi)}{2n} \right].$$

The 95% confidence interval for π is:

$$\left[\frac{zy_1 + y_2}{2n} - 1.96 \frac{\pi(1-\pi)}{2n}, \frac{zy_1 + y_2}{2n} + 1.96 \frac{\pi(1-\pi)}{2n} \right]$$

(Question2): (15 points) Suppose you flip a coin 314 times, and heads appears 159 times.

- Construct an approximate 90% confidence interval for the probability that the coin comes up heads.
- Approximately how many samples would you need to obtain an approximate 90% confidence interval with width 0.02, while keeping exactly 50.64% of flips appearing heads?
- You give the coin to a friend, who also flips the coin 314 times, and obtains an approximate confidence interval [0.390, 0.500] instead. What confidence level did (s)he use?

$$(a) \quad \hat{p} = \frac{x}{n} = \frac{159}{314} = 0.506$$

$$1-\alpha = 90\% \Rightarrow \alpha = 10\% = 0.1$$

$$z_{\frac{\alpha}{2}} = z_{0.05} = 1.645$$

$$\frac{\sigma^2}{n} = \frac{\hat{p}(1-\hat{p})}{n} = \frac{\frac{159}{314}(1-\frac{159}{314})}{314}$$

$$\sqrt{\frac{\sigma^2}{n}} = 0.028$$

$$P(-z_{\frac{\alpha}{2}} < \frac{\hat{p} - p}{\sigma/\sqrt{n}} < z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$\text{The lower confidence limit } \hat{L} = \hat{p} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 0.460$$

$$\text{The upper confidence limit } \hat{U} = \hat{p} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 0.552$$

The confidence interval is $[0.460, 0.552]$

$$(b) \hat{p}' = 50.64\% = 0.5064$$

$$\text{width} = \hat{U}' - \hat{L}' = 2 \cdot z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 0.02 \Rightarrow z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 0.01$$

$$z_{\frac{\alpha}{2}} = 1.645$$

$$\sigma^2 = \hat{p}'(1 - \hat{p}') = 0.250$$

$$\Rightarrow \sigma' = 0.5$$

$$\text{So } 1.645 \times \frac{0.5}{\sqrt{n'}} = 0.01 \Rightarrow \sqrt{n'} = 82.25$$

we get: $n' = 6765$.

We need 6765 samples.

$$(c) \text{ The width is } 0.500 - 0.390 = 0.110$$

Suppose the confidence level is $1 - \alpha'$

The original width of confidence interval is $0.552 - 0.460 = 0.092$

$$\frac{z_{\frac{\alpha'}{2}}}{z_{\frac{\alpha}{2}}} = \frac{0.110}{0.092} \Rightarrow z_{\frac{\alpha'}{2}} = 1.97$$

$$\text{So } \frac{\alpha'}{2} = 0.02442 \Rightarrow \alpha' = 0.04884$$

$$1 - \alpha' = 0.95116 = 95\%$$

(Question3): (15 points) Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu_X, \sigma_X^2)$, $Y_1, \dots, Y_m \stackrel{i.i.d.}{\sim} N(\mu_Y, \sigma_Y^2)$ and they are independent. Suppose that both σ_X^2 and σ_Y^2 are known.

(a) Construct a two-sided $100(1 - \alpha)\%$ confidence interval for the difference $\mu_X - \mu_Y$.

(b) We want to obtain a 90% confidence interval for the difference between true average cable strengths made by Company X and by Company Y. Suppose cable strength is normally distributed for both types of cables with $\mu_X = 50$ and $\mu_Y = 30$. If we can make $n + m = 6000$ observations, how many of these should be on Company X cable if we want to minimize the width of the interval?

$$(a) \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}} \sim N(0, 1),$$

$$\text{where } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i$$

$$P(-z_{\frac{\alpha}{2}} < \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\sigma_X^2/n + \sigma_Y^2/m}} < z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$\Rightarrow \hat{L} = \bar{X} - \bar{Y} - \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \cdot Z_{\frac{\alpha}{2}}$$

$$\hat{U} = \bar{X} - \bar{Y} + \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \cdot Z_{\frac{\alpha}{2}}$$

Two-sided $100(1-\alpha)\%$ confidence interval is:

$$\left[\bar{X} - \bar{Y} - \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \cdot Z_{\frac{\alpha}{2}}, \bar{X} - \bar{Y} + \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \cdot Z_{\frac{\alpha}{2}} \right]$$

$$(b) 1 - \alpha = 90\% \Rightarrow \alpha = 10\% = 0.1$$

The width of the confidence interval is $2Z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}$

$$Z_{\frac{\alpha}{2}} = Z_{0.05} = 1.645.$$

To minimize the width of the CI is to

minimize $\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}$ i.e., minimize $(\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m})$.

Given that $n+m=6000$

$$\text{Let } f(n) = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m} = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{6000-n}$$

$$\begin{aligned} f'(n) &= -\frac{\sigma_X^2}{n^2} + \frac{\sigma_Y^2}{(6000-n)^2} \\ &= \frac{n^2\sigma_Y^2 - (6000-n)^2\sigma_X^2}{n^2(6000-n)^2} \end{aligned}$$

$$\text{Let } f'(n)=0, (\sigma_X^2 - \sigma_Y^2)n^2 - 12000\sigma_X^2 n + 36000000\sigma_X^2 = 0.$$

Solve the equation, we get

$$\begin{aligned} n &= \frac{12000(\sigma_X^2 - \sigma_X\sigma_Y)}{2(\sigma_X^2 - \sigma_Y^2)} \\ &= \frac{6000\sigma_X}{\sigma_X + \sigma_Y} = \frac{6000 \times 50}{50 + 30} = 3750 \end{aligned}$$

(Question4): (15 points) Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$. In this question, we will construct confidence interval for σ^2 . Let $C \sim \chi^2(r)$, and as usual for any $\alpha \in [0, 1]$ we denote $\chi_\alpha^2(r)$ to be

$$P(C > \chi_\alpha^2(r)) = \alpha.$$

(a) Suppose that μ is known. By considering the distribution of

$$\sum_i \left(\frac{X_i - \mu}{\sigma} \right)^2,$$

prove that

$$\left[\frac{\sum_i (X_i - \mu)^2}{\chi_{\alpha/2}^2(n)}, \frac{\sum_i (X_i - \mu)^2}{\chi_{1-\alpha/2}^2(n)} \right]$$

is a $100(1-\alpha)\%$ confidence interval for σ^2 .

Proof: Since μ is known,

$$\text{the sample variance } S^2 = \frac{1}{n} \sum_i (X_i - \mu)^2$$

$$\sum \left(\frac{X_i - \mu}{\sigma} \right)^2 = \frac{1}{\sigma^2} \sum (X_i - \mu)^2 = \frac{nS^2}{\sigma^2} \sim \chi^2(n)$$

$$\text{Since } P(C > \chi^2_{\alpha}(r)) = \alpha,$$

$$P(C > \chi^2_{\frac{\alpha}{2}}(r)) = \frac{\alpha}{2}, \quad P(C > \chi^2_{1-\frac{\alpha}{2}}(r)) = 1 - \frac{\alpha}{2}$$

$$\begin{aligned} P(\chi^2_{1-\frac{\alpha}{2}}(n) \leq \frac{nS^2}{\sigma^2} \leq \chi^2_{\frac{\alpha}{2}}(n)) \\ = P\left(\frac{nS^2}{\sigma^2} > \chi^2_{1-\frac{\alpha}{2}}(n)\right) - P\left(\frac{nS^2}{\sigma^2} > \chi^2_{\frac{\alpha}{2}}(n)\right) \\ = 1 - \alpha. \end{aligned}$$

$$\text{i.e., } P\left(\chi^2_{1-\frac{\alpha}{2}}(n) \leq \frac{\sum_i (X_i - \mu)^2}{\sigma^2} \leq \chi^2_{\frac{\alpha}{2}}(n)\right) = 1 - \alpha$$

$$\Rightarrow P\left(\frac{\sum_i (X_i - \mu)^2}{\chi^2_{\frac{\alpha}{2}}(n)} \leq \sigma^2 \leq \frac{\sum_i (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n)}\right) = 1 - \alpha$$

$$\text{So } \left[\frac{\sum_i (X_i - \mu)^2}{\chi^2_{\frac{\alpha}{2}}(n)}, \frac{\sum_i (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n)} \right] \text{ is a } 100(1-\alpha)\% \text{ confidence interval for } \sigma^2.$$

(b) Suppose that μ is unknown. By considering the distribution of

$$\sum_i \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 = \frac{(n-1)S^2}{\sigma^2},$$

prove that

$$\left[\frac{\sum_i (X_i - \bar{X})^2}{\chi^2_{\alpha/2}(n-1)}, \frac{\sum_i (X_i - \bar{X})^2}{\chi^2_{1-\alpha/2}(n-1)} \right]$$

is a $100(1-\alpha)\%$ confidence interval for σ^2 .

Proof: μ is unknown.

$$\text{Sample variance } S^2 = \frac{1}{n-1} \sum_i \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 = \frac{1}{n-1} \cdot \frac{1}{\sigma^2} \sum (X_i - \bar{X})^2$$

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \quad \text{i.e., } \frac{\sum_i (X_i - \bar{X})^2}{\sigma^2} \sim \chi^2(n-1)$$

$$\text{Since } P(C > \chi^2_{\alpha}(r)) = \alpha,$$

$$P(C > \chi^2_{\frac{\alpha}{2}}(r)) = \frac{\alpha}{2}, \quad P(C > \chi^2_{1-\frac{\alpha}{2}}(r)) = 1 - \frac{\alpha}{2}$$

$$\begin{aligned} P(\chi^2_{1-\frac{\alpha}{2}}(n-1) \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi^2_{\frac{\alpha}{2}}(n-1)) \\ = P\left(\frac{(n-1)S^2}{\sigma^2} > \chi^2_{1-\frac{\alpha}{2}}(n-1)\right) - P\left(\frac{(n-1)S^2}{\sigma^2} > \chi^2_{\frac{\alpha}{2}}(n-1)\right) \\ = 1 - \alpha. \end{aligned}$$

$$\text{i.e., } P\left(\chi^2_{1-\frac{\alpha}{2}}(r) \leq \frac{\sum_i (X_i - \mu)^2}{\sigma^2} \leq \chi^2_{\frac{\alpha}{2}}(r)\right) = 1 - \alpha$$

$$\Rightarrow P\left(\frac{\sum_i (X_i - \mu)^2}{\chi^2_{\frac{\alpha}{2}}(n-1)} \leq \sigma^2 \leq \frac{\sum_i (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)}\right) = 1 - \alpha$$

So $\left[\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)}, \frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)} \right]$ is a $100(1-\alpha)\%$ confidence interval for σ^2 .

- (c) In the same setting as part (b), that is, suppose that μ is unknown. Construct a $100(1-\alpha)\%$ confidence interval for σ .

Know that $\left[\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)}, \frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)} \right]$ is a $100(1-\alpha)\%$ confidence interval for σ^2 .

$$\sigma = \sqrt{\sigma^2} > 0$$

\Rightarrow The $100(1-\alpha)\%$ confidence interval for σ

$$\text{is: } \left[\sqrt{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)}}, \sqrt{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi^2_{1-\frac{\alpha}{2}}(n-1)}} \right]$$

In [8]:

```

# STA2004 Programming assignment2
# Name: Ou Ziyi Student ID:121090429

# 2.1
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
from functools import reduce

## Since  $\mu = E[Y] = 1/\lambda$  and  $y\_bar = 10$ , we can know that  $\lambda\_hat = 0.1$ 
##  $L(\lambda) = \lambda^2 * e^{-\lambda ny}$ 

fig, ((ax1, ax2, ax3, ax4)) = plt.subplots(4, 1, figsize=(8, 14), sharex=True)
n_value = np.arange(0.01, 1.00, 0.01)
list1= []
list2 =[]
for i in n_value:
    b = pow(i, 1)*np.exp(-i*10*1) - pow(1/10, 1)*np.exp(-1/10*10*1)
    list1.append(i)
    list2.append(b)
ax1.set_title('n = {}'.format(1))
ax1.plot(list1, list2)

list1= []
list2 =[]
for i in n_value:
    b = pow(i, 5)*np.exp(-i*10*5) - pow(1/10, 5)*np.exp(-1/10*10*5)
    list1.append(i)
    list2.append(b)
ax2.set_title('n = {}'.format(5))
ax2.plot(list1, list2, color = 'green')

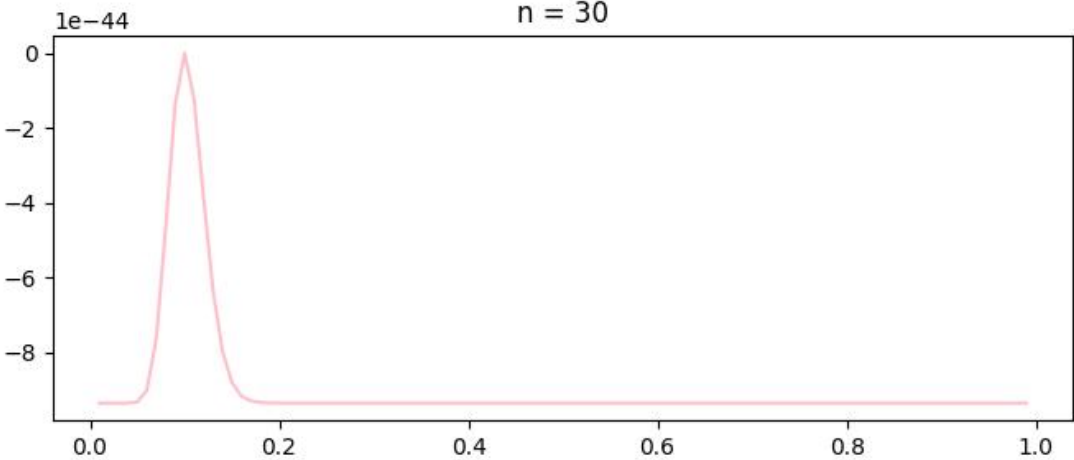
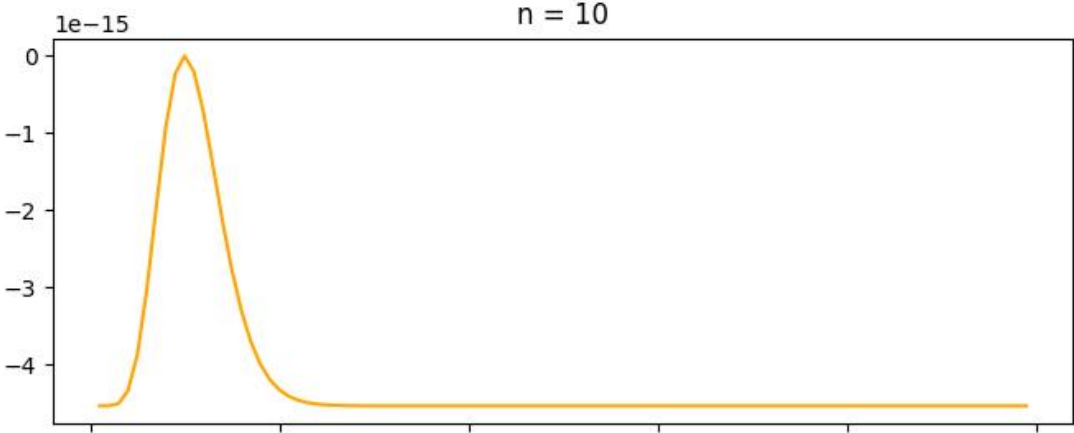
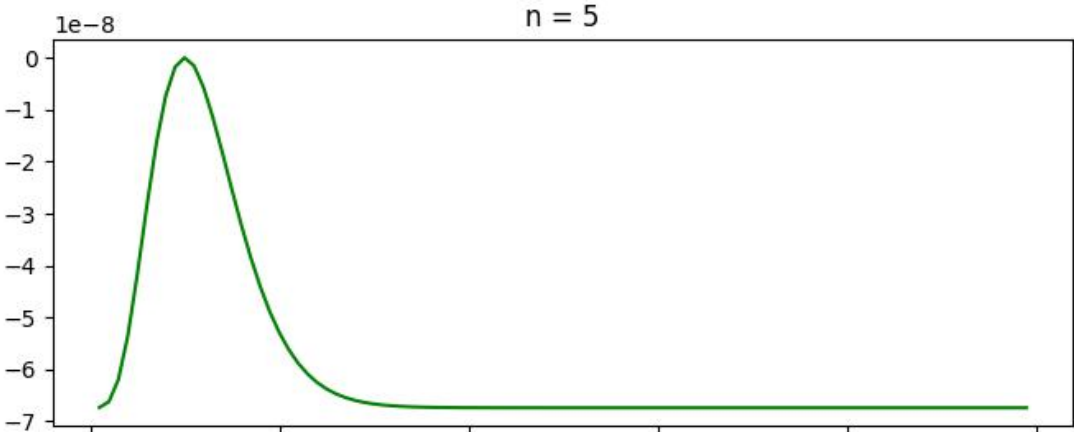
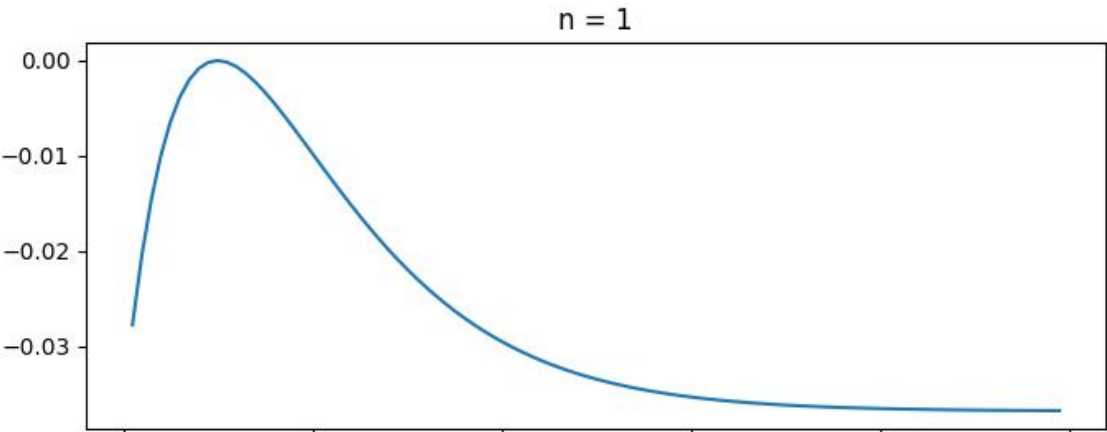
list1= []
list2 =[]
for i in n_value:
    b = pow(i, 10)*np.exp(-i*10*10) - pow(1/10, 10)*np.exp(-1/10*10*10)
    list1.append(i)
    list2.append(b)
ax3.set_title('n = {}'.format(10))
ax3.plot(list1, list2, color = 'orange')

list1= []
list2 =[]
for i in n_value:
    b = pow(i, 30)*np.exp(-i*10*30) - pow(1/10, 30)*np.exp(-1/10*10*30)
    list1.append(i)
    list2.append(b)
ax4.set_title('n = {}'.format(30))
ax4.plot(list1, list2, color = 'pink')

```

Out[8]:

[<matplotlib.lines.Line2D at 0x21fc7fb8880>]



In [7]:

```
# 2.2
import numpy as np
n_value = [1, 5, 10, 30]
for n in n_value:
    confidence_interval = np.array([1/10 - (1/10*1.96)/np.sqrt(n), 1/10 + (1/10*1.96)/np.sqrt(n)]
    print(confidence_interval)
```

```
[-0.096  0.296]
[0.01234614 0.18765386]
[0.03801936 0.16198064]
[0.06421546 0.13578454]
```

In [17]:

```

# 2.3
lambda_zero = 1
Z_half_alpha = 1.96
size = 10000
plt.title('Confidence Intervals')
plt.xlabel('lambda')
plt.ylabel('Samples')
plt.vlines(1, 0, 100, colors='green')
for num in range(100):
    data_set = np.random.exponential(lambda_zero, size)
    mean_of_data = np.mean(data_set)
    lambda_hat = 1/mean_of_data
    half_length = lambda_hat * Z_half_alpha/np.sqrt(size)
    if lambda_hat - half_length > 1 or lambda_hat + half_length < 1:
        plt.hlines(num, lambda_hat - half_length, lambda_hat + half_length, colors='red')
    else:
        plt.hlines(num, lambda_hat - half_length, lambda_hat + half_length, colors='black')
plt.show()

```

