

CS 237 Fall 2019 Homework Seven

Due date: PDF file due Friday October 25th @ 11:59PM in GradeScope with 12-hour grace period

No late submissions!

General Instructions

Please complete this notebook by filling in solutions where indicated. Be sure to "Run All" from the Cell menu before submitting.

There are two sections to the homework: problems 1 - 10 are analytical problems about last week's material, and the remaining problems are coding problems which will be discussed in lab next week.

In [1]:

```

# Useful imports and definitions for CS 237

import numpy as np                # arrays and functions which operate on array
from numpy import linspace, arange, mean
import matplotlib.pyplot as plt   # normal plotting
import seaborn as sns             # Fancy plotting
import pandas as pd              # Data input and manipulation

from random import random, randint, uniform, choice, sample, shuffle, seed
from collections import Counter
from math import log, floor, ceil, exp

%matplotlib inline

# Calculating permutations and combinations efficiently

def P(N,K):
    res = 1
    for i in range(K):
        res *= N
        N = N - 1
    return res

def C(N,K):
    if(K < N/2):
        K = N-K
    X = [1]*(K+1)
    for row in range(1,N-K+1):
        X[row] *= 2
        for col in range(row+1,K+1):
            X[col] = X[col]+X[col-1]
    return X[K]

def round4(x):
    return round(x+0.00000000001,4)

def round4_list(L):
    return [ round4(x) for x in L]

# Useful code from HW 01

# This function takes a list of outcomes and a list of probabilities and
# draws a chart of the probability distribution.
# It allows labels for x axis with numbers or strings; for the latter, you
# still need to give the numeric labels, but can overwrite them with your string labels

def draw_distribution(Rx, fx, title='Probability Distribution', my_xticks = [], f_size=10):
    plt.figure(figsize=f_size)
    plt.bar(Rx,fx,width=1.0,edgecolor='black')
    plt.ylabel("Probability")
    plt.xlabel("Outcomes")
    if my_xticks != []:
        plt.xticks(Rx, my_xticks)
    elif (Rx[-1] - Rx[0] < 30):
        ticks = range(Rx[0],Rx[-1]+1)
        plt.xticks(ticks, ticks)
    plt.title(title)

```

```

plt.show()

# Example of use

# draw_distribution([1,2,3,4], [0.25,0.35,0.15,0.25])

# p = 0.14159234368

# draw_distribution([0,1], [p,1.0-p], "Distribution for Unfair Coin", ['Heads', 'Ta

# This draws a useful bar chart for the distribution of the
# list of integer in outcomes

def show_distribution(outcomes, title='Probability Distribution', my_xticks = [], f
    plt.figure(figsize=f_size)
    num_trials = len(outcomes)
    X = range(int(min(outcomes)), int(max(outcomes))+1)
    freqs = Counter(outcomes)
    Y = [freqs[i]/num_trials for i in X]
    plt.bar(X,Y,width=1.0,edgecolor='black')
    if my_xticks != []:
        plt.xticks(X, my_xticks)
    elif (X[-1] - X[0] < 30):
        ticks = range(X[0],X[-1]+1)
        plt.xticks(ticks, ticks)
    plt.xlabel("Outcomes")
    plt.ylabel("Probability")
    plt.title(title)
    plt.show()

# Example of use

#show_distribution([1,4,3,5,4,6,2,4,3,5,4])

# Scipy statistical functions

from scipy.stats import norm, binom

# https://docs.scipy.org/doc/scipy/reference/stats.html

#### Normal Distribution #####

# Note that in this library loc = mean and scale = standard deviation

# Examples assume random variable X (e.g., housing prices) normally distributed with

# probability density function

norm.pdf(x=50,loc=60, scale= 10)

#a. Find P(X<50)
norm.cdf(x=50,loc=60,scale=10) # 0.4012936743170763

#b. Find P(X>50)
norm.sf(x=50,loc=60,scale=40) # 0.5987063256829237

#c. Find P(60<X<80)
norm.cdf(x=80,loc=60,scale=40) - norm.cdf(x=60,loc=60,scale=40)

#d. What is the minimum cost of the 5% most expensive houses? Alternately, what is

```

```

norm.isf(q=0.05, loc=60, scale=40)

#e. What is the maximum cost of the 5% cheapest houses? Alternately, what is x where
norm.ppf(q=0.05, loc=60, scale=40)

#f give the endpoints of the range for the central alpha percent of the distribution
norm.interval(alpha=0.3, loc=60, scale=140)

#g. generate a random variate
norm.rvs(loc=60, scale=40)

#h. generate random variates, returns list of length = size
norm.rvs(loc=60, scale=40, size=10)

##### Binomial Distribution  $X \sim B(n, p)$  #####

# n = number of independent Bernoulli trials
# p = probability of success for Bernoulli trial
# k = outcome in range [0 .. n]

# Generate a random variate
binom.rvs(n=10, p=0.5)

# Generate a list of random variates
binom.rvs(n=10, p=0.5, size=100)

# Probability mass/density function.
binom.pmf(k=4, n=10, p=0.5)

# Cumulative distribution function
binom.cdf(k=4, n=10, p=0.5)

print()

```

Problem One (Poisson)

This problem is a collection of several different applications of the basic formulae for the Poisson Distribution. You **must** use the Poisson and show all work.

- Assume that the probability that a poker hand is a full house is 0.0014. What is the probability that in 500 random poker hands there are at least two full houses?
- On average, there are three misprints in every 10 pages of a particular book. If every chapter of the book contains 35 pages, what is the probability that Chapters 1 and 5 have exactly 10 misprints each?
- Suppose that on a summer evening, shooting stars are observed at a Poisson rate of one every 12 minutes. What is the probability that three shooting stars are observed in 30 minutes?
- Suppose that in Japan earthquakes occur at a Poisson rate of three per week. What is the probability that the next earthquake occurs after two weeks?

(e) Suppose that, for a telephone subscriber, the number of wrong numbers is Poisson, at a rate of $\lambda = 1$ per week. A certain subscriber has not received any wrong numbers from Sunday through Friday. What is the probability that he receives no wrong numbers on Saturday either?

Hint: Remember that the Poisson Distribution does not have the memory-less property, but each arrival is *independent* of every other arrival, and in particular, two intervals of time or space which do not overlap have an independent number of arrivals; so, like coins, arrivals don't have memory...

Solution

- (a) Expected Value = $0.0014 * 500 = 0.7$
 $P(\text{at least two full house out of 500})$
 $= 1 - P(\text{no full house}) - P(\text{one full house})$
 $= 1 - e^{(-0.7)} * (0.7)^0 / 0! - e^{(-0.7)} * (0.7)^1 / 1!$
 $= 1 - e^{(-0.7)} - e^{(-0.7)} * 0.7$
 $= 1 - 0.4966 - 0.3476$
 $= 0.1558$
- (b) Expected Value = $0.3 * 35 = 10.5$
 $P(\text{Chapter 1 and 5 have 10 misprints})$
 $= P(\text{Chapter 1 has 10 misprints}) * P(\text{Chapter 5 has 10 misprints})$
 $= e^{(-10.5)} * (10.5)^{10} / 10! + e^{(-10.5)} * (10.5)^{10} / 10!$
 $= 0.1236 + 0.1236$
 $= 0.2472$
- (c) Expected Value = $1/12 * 30 = 2.5$
 $P(\text{three shooting stars observed})$
 $= e^{(-2.5)} * (2.5)^3 / 3!$
 $= 0.2138$
- (d) Expected Value = $3 * 2 = 6$
 $P(\text{no earthquake for 2 weeks})$
 $= e^{(-6)} * (6)^0 / 0!$
 $= e^{(-6)}$
 $= 0.002479$
- (e) Expected Value = $1/7 * 1 = 1/7$
 $P(\text{no wrong numbers on Saturday either})$
 $= e^{(-1/7)} * (1/7)^0 / 0!$
 $= 0.8669$

Problem Two (Poisson compared with Exponential)

Suppose the emergency room at Mass General opens at 6am and has a mean arrival rate throughout the day of $\lambda = 6.9$ patients per hour.

You must use the Poisson for (a) - (c) and Exponential for (d) and (e). Give the name of the distribution and show the parameters (for example, (a) is $P(X=12)$ for $X \sim \text{Poi}(6.9)$), and show all work.

- (a) What is the probability that exactly 12 patients arrive between 6am and 7am? (Use Poisson)
- (b) What is the probability that no patient arrives before 7am? (Use Poisson)
- (c) What is the probability that a patient arrives between 6:15 and 6:45? (Use Poisson)
- (d) What is the probability that the first patient arrives between 6am and 7am? (Use Exponential)

(e) What is the probability that no patient arrives between 6:15 and 6:45? (Use Exponential)

Hint: Some of these are complements of each other...

Solution:

- (a) $P(X=12)$ for $X \sim \text{Poi}(6.9)$
 $= e^{-6.9} * (6.9)^{12}/12!$
 $= 0.02450$
- (b) $P(\text{no patient arrives before 7am})$
 $= P(X=0)$ for $X \sim \text{Poi}(6.9)$
 $= e^{-6.9} * (6.9)^0/0!$
 $= 0.001008$
- (c) $P(\text{patient arrives between 6:15 and 6:45})$
 $= P(X=1)$ for $X \sim \text{Poi}(6.9)$
 $= e^{-6.9} * (6.9)^1/1!$
 $= 0.006954$
- (d) $P(\text{first patient arrives between 6am and 7am})$
 $= P(1 < X)$ for $X \sim \text{Exp}(6.9)$
 $= 1 - e^{-6.9}$
 $= 0.9990$
- (e) $P(\text{no patient arrives between 6:15 and 6:45})$
 $= P(0 < X)$ for $X \sim \text{Exp}(6.9)$
 $= 1 - e^{-3.45}$
 $= 0.9683$

Problem Three (Exponential)

In the following assume that we are dealing with Poisson processes and can use the Exponential distribution.

(a) Suppose that every three months, on average, an earthquake occurs in California. What is the probability that the next earthquake occurs after three but before seven months?

(b) Suppose we model time to failure of TV tubes as an Exponential random variable and that tubes fail on average after 10 years. If Jim bought his TV set 10 years ago, what is the probability that its tube will last another 10 years?

(c) Guests arrive at a hotel at an average rate of 5 per hour. Suppose that for the last 10 minutes no guest has arrived. What is the probability that the next one will arrive in less than 2 minutes?

(d) Considering the same situation as in (c), what is the probability that from the arrival of the tenth to the arrival of the eleventh guest takes no more than 2 minutes?

Hint: Make sure you understand the difference between the rate parameter λ (= mean number of arrivals per unit time) and $\beta = 1/\lambda$ (= mean interarrival time).

Solution:

- (a) $P(3 < X < 7)$ for $X \sim \text{Exp}(1/3)$
 $= (1 - e^{(-7/3)}) - (1 - e^{(-3/3)})$
 $= 0.9030 - 0.6321$
 $= 0.2709$
- (b) $P(X > 10)$ for $X \sim \text{Exp}(10)$
 $= 1 - P(X < 10)$ for $X \sim \text{Exp}(10)$
 $= 1 - (1 - e^{(-1)})$
 $= e^{(-1)}$
 $= 0.3679$
- (c) $P(X < 2)$ for $X \sim \text{Exp}(12)$
 $= 1 - e^{(-2/12)}$
 $= 0.1535$
- (d) $P(X < 2)$ for $X \sim \text{Exp}(12)$
 $= 1 - e^{(-2/12)}$
 $= 0.1535$

Problem Four

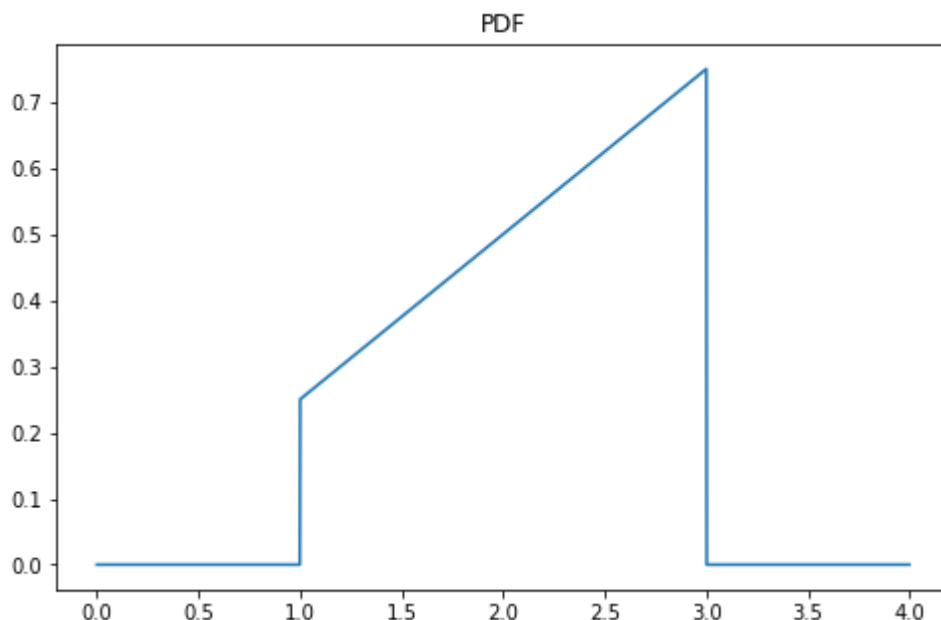
(Continuous Distributions) Let X be a continuous random variable with a frequency distribution (PMF) of the form

$$f(x) = \begin{cases} \frac{x}{4} & \text{if } 1 \leq x \leq 3 \\ 0 & \text{otherwise} \end{cases}$$

which can be graphed as follows:

In [2]:

```
plt.figure(figsize=(8, 5))
plt.title("PDF")
plt.plot(np.arange(0,4,0.001),[x/4 if 1 <= x <= 3 else 0 for x in np.arange(0,4,0.001)])
plt.show()
```



- (a) Determine the formula for the CDF F_X using geometrical techniques (i.e., not using integrals, but considering what happens to the area to the left of a point a by considering the area of geometrical shapes)
- (b) Determine the formula for the CDF $F_X(x)$ using an integral.
- (c) Plot the CDF $F_X(x)$ (using the code above as a model).
- (d) Find $P(X \geq 2)$
- (e) Find $E(X)$

For (d) -- (e) you must use mathematical techniques and not just calculate it using iterative techniques in Python code. You must show all your work, except that you may use Python to calculate results of mathematical formulae.

Solution: (a) $R_x = \{0, 0.5, 1.0, 1.5, 2, 2.5, 3, 3.5, \dots\}$ $F_x = \{0, 0, 0, 0.15625, 0.375, 0.65625, 1, 1, \dots\}$

(b) Integral $x/4 = x^2/8 + C$ When $x = 3$, $x^2/8 + C = 1$ $9/8 + C = 1$ $C = -1/8$

$$F(x) = \begin{cases} \frac{x^2-1}{8} & \text{if } 1 \leq x \leq 3 \\ 1 & \text{if } x > 3 \\ 0 & \text{if } x < 1 \end{cases}$$

(d) $P(X \geq 2)$

Integral of $x/4$ from 2 to 3
 $= x^2/8$ from 2 to 3
 $= 9/8 - 4/8$
 $= 5/8$

(e) $E(x)$

Integral of $x*x/4$ from 1 to 3
 $=$ Integral of $x^2/4$ from 1 to 3
 $= x^3/12$ from 1 to 3
 $= 27/12 - 1/12$
 $= 26/12$
 $= 13/6$

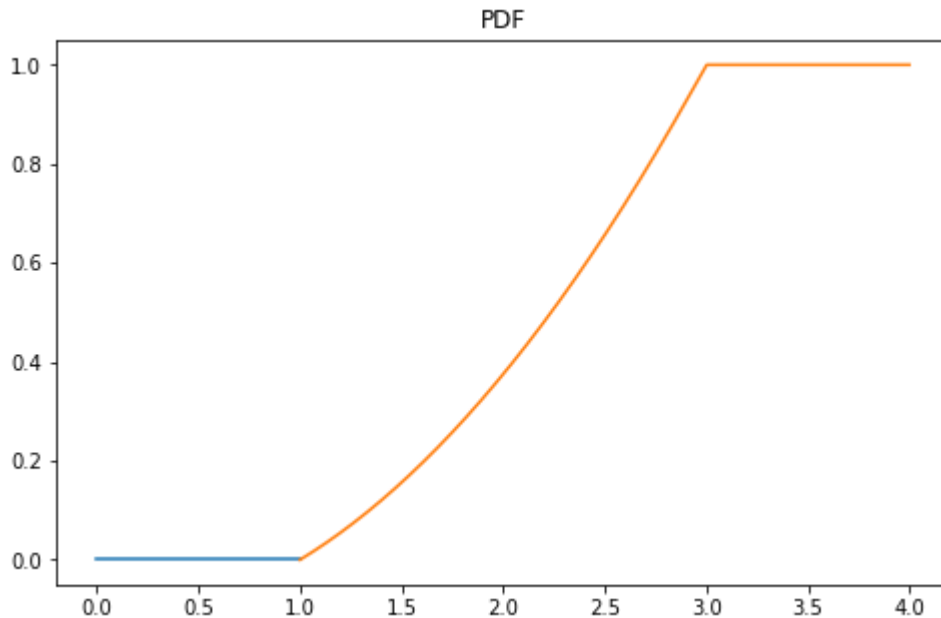
Part c is below...

In []:

In [3]:

```
print("Part c")
plt.figure(figsize=(8, 5))
plt.title("PDF")
plt.plot(np.arange(0,1,0.001),[0 if x <= 1 else 0 for x in np.arange(0,1,0.001)])
plt.plot(np.arange(1,4,0.001),[x*x/8 - 1/8 if 1 <= x <= 3 else 1 for x in np.arange(1,4,0.001)])
plt.show()
```

Part c



Scipy.stats

For the following, use the statistical functions given at the top of this notebook.

Also consider using the Distributions Notebook posted online to visualize these distributions (but use `scipy.stats` for the calculations).

You are not required to do so, but a nice touch is to print out your answer in a code block, i.e., if you were asked "What is the probability in the standard normal that a value occurs in the interval between -0.94 and 1.2 standard deviations from 0," you could answer as follows:

In [4]:

```
mu = 0
sigma = 1
lo = -0.94
hi = 1.2

answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution: " + str(round4(answer)))
```

Solution: 0.7113

Problem Five (Normal Distribution)

Suppose that in a population of individuals, their height is normally distributed with a mean of 68 inches and a standard deviation of 1.45 inches.

- (a) What is the probability that a randomly-selected individual has a height less than 66 inches?
- (b) What is the probability that a randomly-selected individual has a height more than 72 inches?
- (c) What is the probability that a randomly-selected individual has a height between 66.5 and 71 inches?

Solution:

In [5]:

```
print("Part a")
mu = 68
sigma = 1.45
lo = -999999
hi = 66
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution for part a: " + str(round4(answer)))

print()

print("Part b")
lo = 72
hi = 999999
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution for part b: " + str(round4(answer)))

print()

print("Part c")
lo = 66.5
hi = 71
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution for part c: " + str(round4(answer)))
```

Part a
Solution for part a: 0.0839

Part b
Solution for part b: 0.0029

Part c
Solution for part c: 0.8303

Problem Six (Standard Normal Distribution)

Using the same mean and standard deviation as in the previous problem (mean of 68 inches and a standard deviation of 1.45 inches), show how to answer the following questions by first converting to the standard normal distribution (with mean 0 and standard deviation 1), and then calculating the answer in terms of the standard normal (i.e., when you calculate the probabilities, you will be using $N(0,1)$.) Show all work. For (c) you will not need to convert, but simply use the standard normal directly.

- (a) What is the probability that a randomly-selected individual has a height more than 69 inches?
- (b) What is the probability that a randomly-selected individual has a height between 63.2 and 70.7 inches?
- (c) What is the probability that someone is more than 2.7 standard deviations taller than the average height?

Solution:

```
(a) z = (69-68)/1.45
      = 0.69
      P(z > 0.69)
      = 1 - 0.7549
      = 0.2451
(b) z1 = (63.2-68)/1.45
      = -3.31
      z2 = (70.7-68)/1.45
      = 1.86
      P(-3.31 < z < 1.86)
      = 0.9686 - 0.0005
      = 0.9681
(c) P(z > 2.7)
      = 1 - 0.9965
      = 0.0035
```

Problem Seven (Normal Distribution)

Suppose the grades on the 237 midterm are normally distributed with mean 83 and variance of 23. Answer the following questions. By the " k^{th} percentile" we mean the x value for which $P(X \leq x) = k/100$ (meaning k is in percentages, and $k/100$ is in terms of probabilities.)

Browse the functions from `scipy.stats.norm` given to see which one will solve each problem.

- (a) If an A is 93 or above, what percentage of the class will get an A?
- (b) What is the exam score which represents the 90th percentile? (Assume exam scores are real numbers.)
- (c) If a score between 70 and 80 results in a C, and there are 120 people in the class, approximately how many people will get a C? (Round to the nearest integer.)

Solution:

In [6]:

```
print("Part a")
mu = 83
sigma = 23 ** 0.5
lo = 93
hi = 999999
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution for part a: " + str(round4(answer)))

print()

print("Part b")
percent = norm.ppf(q = 0.9, loc = mu, scale = sigma)
print("Exam score which represents 90th percentile is " + str(round4(percent)))

print()

print("Part c")
mu = 83
sigma = 23 ** 0.5
lo = 70
hi = 80
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Solution for part c: " + str(round4(answer)))
```

Part a

Solution for part a: 0.0185

Part b

Exam score which represents 90th percentile is 89.1461

Part c

Solution for part c: 0.2625

Problem Eight (Combining Poisson with other Distributions)

Suppose that in a certain region of California, earthquakes occur at the average rate of 7 per year.

- (a) What is the probability that in exactly three of the next eight years, no earthquakes occur?
- (b) What is the **expect**ed number of years to wait until we have a year with exactly 7 earthquakes?
- (c) In the next century, how many years would you **expect** to see with more than 10 earthquakes?

Hint: When you see the word "expect" you should expect to use the expected value!

Solution

(a) Expected Value = 7

$$\begin{aligned} &P(\text{exactly three of the next eight years, no earthquake occur}) \\ &= C(8,3) * (e^{-7}) * 7^3 / 3! * (1 - (e^{-7}) * 7^0 / 0!)^5 \\ &= 4.227 * 10^{-8} \text{ or } 0.00000004227 \end{aligned}$$

(b) Expected number of years to wait until we have a year with exactly 7 earthquakes

$$\begin{aligned} &= 1/P(\text{exactly 7 earthquakes}) \\ &= 1/(e^{-7}) * 7^7 / 7! \\ &= 1/0.1490 \\ &= 6.7114 \end{aligned}$$

The expected number of years to wait until we have a year with exactly 7 earthquakes is 6.7114 years.

(c) Expected years in next century with more than 10 earthquakes

$$\begin{aligned} &= 100 * P(11 \text{ or more earthquakes occurring}) \\ &= 100 * (1 - P(10 \text{ or less earthquakes occurring})) \\ &= 100 * (1 - (e^{-7}) * 7^0 / 0! + e^{-7}) * 7^1 / 1! + e^{-7}) * 7^2 / 2! \\ &+ \dots + e^{-7}) * 7^9 / 9! + e^{-7}) * 7^{10} / 10!)) \\ &= 100 * (1 - 0.9015) \\ &= 100 * 0.09852 \\ &= 9.852 \end{aligned}$$

You would expect 9.852 years to see with more than 10 earthquakes

Problem Nine (Combining Normal with Other Distributions)

Suppose that in the Men's Olympic Ski Team, the chest size measurements are normally distributed with a mean of 39.8 inches and a standard deviation of 2.05 inches.

(a) What the probability that of 20 randomly selected members of the team, at least 5 have a chest size of at least 41.7 inches?

(b) Supposing we choose men on the team repeatedly and with replacement, how many men would you expect to choose before finding a member with a chest measurement of less than 37 inches?

Hint: Let X is a normally distributed random variable according to the parameters given in the first sentence. Then consider Y and Z be appropriately distributed random variables for (a) and (b) respectively.

Solution

In [7]:

```

print("Part a")
mu = 39.8
sigma = 2.05
lo = 41.7
hi = 9999999
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
final_answer = C(20,0) * answer ** 0 * (1-answer) ** 20 + C(20,1) * answer ** 1 * (1-answer) ** 19 +
C(20,2) * answer ** 2 * (1-answer) ** 18 + C(20,3) * answer ** 3 * (1-answer) ** 17 +
C(20,4) * answer ** 4 * (1-answer) ** 16
print("Solution for part a: " + str(round4(final_answer)))

print()

print("Part b")
mu = 39.8
sigma = 2.05
lo = -9999999
hi = 37
answer = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
final_answer = 1/answer
print("You would expect to choose " + str(round4(final_answer)) +
      " men to find a member with a chest measurement of less than 37 inches")

```

Part a

Solution for part a: 0.7275

Part b

You would expect to choose 11.6289 men to find a member with a chest measurement of less than 37 inches

Problem Ten (Normal Approximation to the Binomial)

This problem concerns the normal approximation to the binomial. The "continuity correction" (sometimes called Yates's Continuity Correction) is a technique I will show in the Monday 10/21 lecture for improving the accuracy of the normal approximation.

In this problem we will measure the accuracy of the approximations by using **absolute percentage error**, which is defined to be:

$$\frac{|\text{measured value} - \text{actual value}|}{\text{actual value}} \times 100.$$

(a) Suppose of all the kids that show up on Halloween night, 58% are dressed in Deadpool costumes. If 60 kids show up, what is the probability that between 33 and 38 (inclusive) kids will be dressed in Deadpool costumes? (Use the binomial.)

(b) Repeat the previous question, but using the normal approximation to the binomial, without using the continuity correction, and express the accuracy of your approximation using the absolute percentage error.

(c) Repeat the previous question, but now using the continuity correction, again showing the accuracy using the absolute percentage error.

Solution:

In [8]:

```

print("Part a")
answer1 = C(60,33) * (0.58) ** 33 * (0.42) ** 27 + C(60,34) * (0.58) ** 34 * (0.42) ** 26
print("Probability that between 33 and 38 (inclusive) kids will be dressed in Deadpool costumes is", answer1)

print()
print("Part b")
mu = 60 * 0.58
sigma = (60 * 0.58 * 0.42) ** 0.5
lo = 33
hi = 38
answer2 = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Normal Approximation to the binomial for part a without using continuity correction is", answer2)

value = abs(((answer2) - (answer1))) / answer1 * 100
print("The absolute percentage error is " + str(round4(value)) + "%.")

print()
print("Part c")
mu = 60 * 0.58
sigma = (60 * 0.58 * 0.42) ** 0.5
lo = 32.5
hi = 38.5
answer2 = norm.cdf(x=hi,loc=mu,scale=sigma) - norm.cdf(x=lo,loc=mu,scale=sigma)
print("Normal Approximation to the binomial for part a using continuity correction is", answer2)

value = abs(((answer2) - (answer1))) / answer1 * 100
print("The absolute percentage error is " + str(round4(value)) + "%.")

```

Part a

Probability that between 33 and 38 (inclusive) kids will be dressed in Deadpool costumes is 0.5609

Part b

Normal Approximation to the binomial for part a without using continuity correction is 0.4798

The absolute percentage error is 14.4532%.

Part c

Normal Approximation to the binomial for part a using continuity correction is 0.5597

The absolute percentage error is 0.2107%.

Lab Problems

Generating Variates from a Continuous Distribution

Samples from a given distribution are often called "random variates" or just "variates" for short; to generate **size** random variates from a normal distribution with mean **loc** and standard deviation **scale** we can use the `scipy.stats` function

```
X = norm.rvs(loc=0,scale=1,size=num_trials)
```

(Note that in this case, the normal is defined in terms of the standard deviation, and **not** the variance.)

Run the next cell several times to get a sense for how this function works

In [9]:

```
print("From N(0,1):")
X = norm.rvs()      # default is a standard normal with mean 0 and standard deviation 1
print(X)
print()
print("From N(10,2^2):")
X = norm.rvs(10,2)   # defined by mean and standard deviation (NOT the variance)
print(X)
print()
print("Ten variates from N(66,3^2):")
X = norm.rvs(66,3,size=10)
print(X)
```

```
From N(0,1):
-1.4359925963667715
```

```
From N(10,2^2):
14.489388681989961
```

```
Ten variates from N(66,3^2):
[61.81835132 67.25593124 68.30547861 61.64468197 68.64152813 61.969505
 2
 63.22132448 64.9914444 68.81437638 64.08049588]
```

Graphing Continuous Variates: A Problem

So generating normal variates is easy! Think of them as heights of BU students picked at random. We'll assume they have a mean height of 66 inches with a variance of 3 inches.

What we are going to concern ourselves with in this first problem is now to display a collection of such normal numbers.

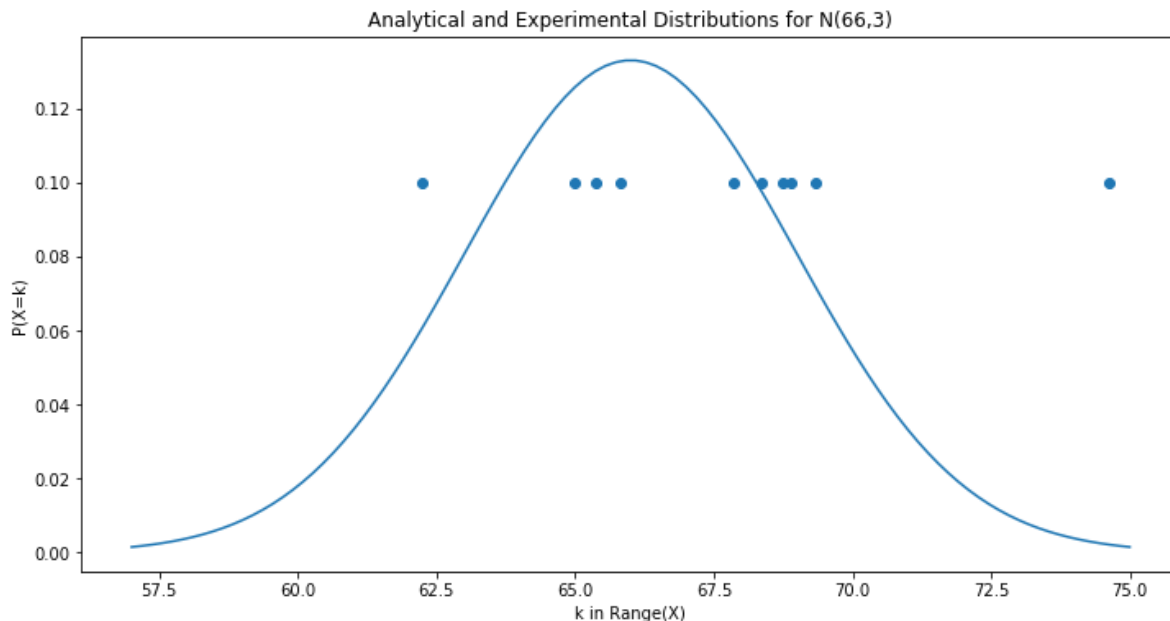
Here is the problem: since each value occurs (with high probability) only once, we can't just create a histogram and convert it into a frequency distribution. **Essentially we are trying to create a frequency distribution from values whose theoretical probability is 0.**

Here is what happens if we do this, and graph it as a scatter plot against the theoretical (continuous) distribution:

In [10]:

```
def display_normal_samples(mu,sigma,num_trials):
    fig, ax = plt.subplots(1,1,figsize=(12,6))
    plt.title('Analytical and Experimental Distributions for N('+str(mu)+'','+str(sigma))')
    plt.ylabel("P(X=k)")
    plt.xlabel("k in Range(X)")
    # use normal(...) to generate random samples
    X = sorted(norm.rvs(mu,sigma,num_trials))
    # Now convert frequency counts into probabilities
    D = Counter( X )
    P = [D[k]/num_trials for k in X]
    plt.scatter(X,P)
    # Now generate the theoretical normal with the same mean and
    X2 = np.linspace(mu-sigma*3,mu+sigma*3,100)
    Y = [norm.pdf(x,mu,sigma) for x in X2]
    plt.plot(X2,Y)
    plt.show()

# try setting the number of trials - the number of samples generated -- to 100 and
num_trials = 10
display_normal_samples(66,3,num_trials)
```



Graphing Continuous Variates with Bins

You see the problem: since each floating point number (an approximation of a real number) is generated with high probability at most once, we can't see the accumulation of samples that would indicate the probability. What to do? Well, you know that probabilities can only be calculated in continuous distributions using **intervals**, so we will create intervals by **rounding** our floating-point samples to put the variates into various-width bins, and then plot the probabilities of each bin.

So we must specify the width of each bin that we want to use as an interval to collect together our samples from the continuous distribution. Then we must "slot" each variate into its appropriate bin.

Graphing Normal Variates with Bins Calculated from Standard Deviation

For the normal distribution it makes sense to define the bin boundaries in terms of standard deviations from the mean, since we will be dealing with an unknown range of data; this bins can be made as wide or as narrow as we want, but will represent an interval defined in terms of the standard deviation sigma of the distribution.

We will graph the distribution in a range of at least 4 standard deviations of the mean, ignoring the rare occurrence of a variate outside this range.

In [11]:

```
# Define the boundaries of bins with the specified width around the mean,
# to plus/minus at least 4 * sigma

# bin_width is in units of sigma, so bin_width = 0.1 means sigma/10

def makeBins(mu,sigma,bin_width):
    numBins = ceil(4/bin_width)
    bins = [mu+sigma*bin_width*x for x in range(-numBins,numBins+1)]
    return bins

# Change the parameters several times to see the effect of this

print(makeBins(0,1,.5))
```

```
[-4.0, -3.5, -3.0, -2.5, -2.0, -1.5, -1.0, -0.5, 0.0, 0.5, 1.0, 1.5,
 2.0, 2.5, 3.0, 3.5, 4.0]
```

Problem Eleven: Generating and Graphing Normal Variates

How does the number of trials affect the fit of the data to the normal distribution?

Now let's do our previous experiment but trying various values for bin_width...

In [12]:

```
def display_normal_samples_binned(mu,sigma,num_trials,bin_widths):
    fig, ax = plt.subplots(1,1,figsize=(12,6))
    plt.title('Analytical and Experimental Distributions for N('+str(mu)+'',' +str(sigma)')
    plt.ylabel("f(k)")
    plt.xlabel("k in Range(X)")
    # use norm.rvs(...) to generate random samples
    X = norm.rvs(mu,sigma,num_trials)
    plt.hist(X,bins=makeBins(mu,sigma,bin_widths),normed=True,edgecolor='k',alpha=0.5)
    # if get deprecation warning for normed, use the next line instead (or don't worry)
    # plt.hist(X,bins=makeBins(mu,sigma,bin_widths),density=True,edgecolor='k',alpha=0.5)

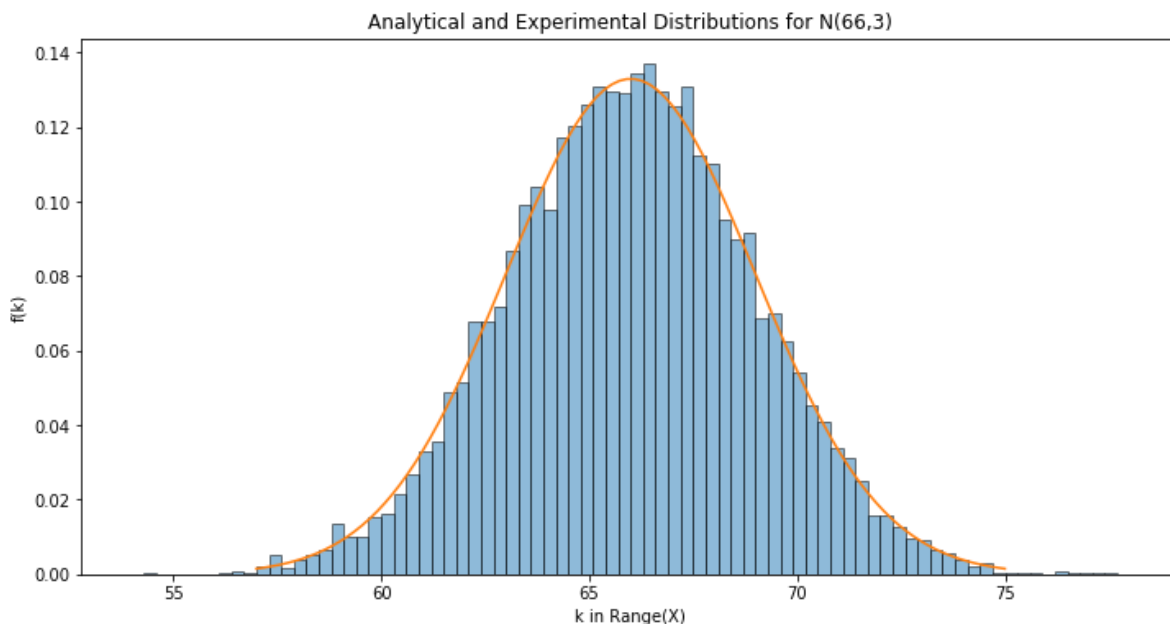
    # Now generate the theoretical normal with the same mean and
    X2 = np.linspace(mu-sigma*3,mu+sigma*3,100)
    Y = [norm.pdf(x,mu,sigma) for x in X2]
    plt.plot(X2,Y)
    plt.show()

#try each of these and observe the effects

N = 10000    # try 100, 1000, and 1,000,000

display_normal_samples_binned(66,3,N,0.1)
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:8: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.



Clearly the data seems to fit the normal better when the number of trials increases...

Affect of the bin width

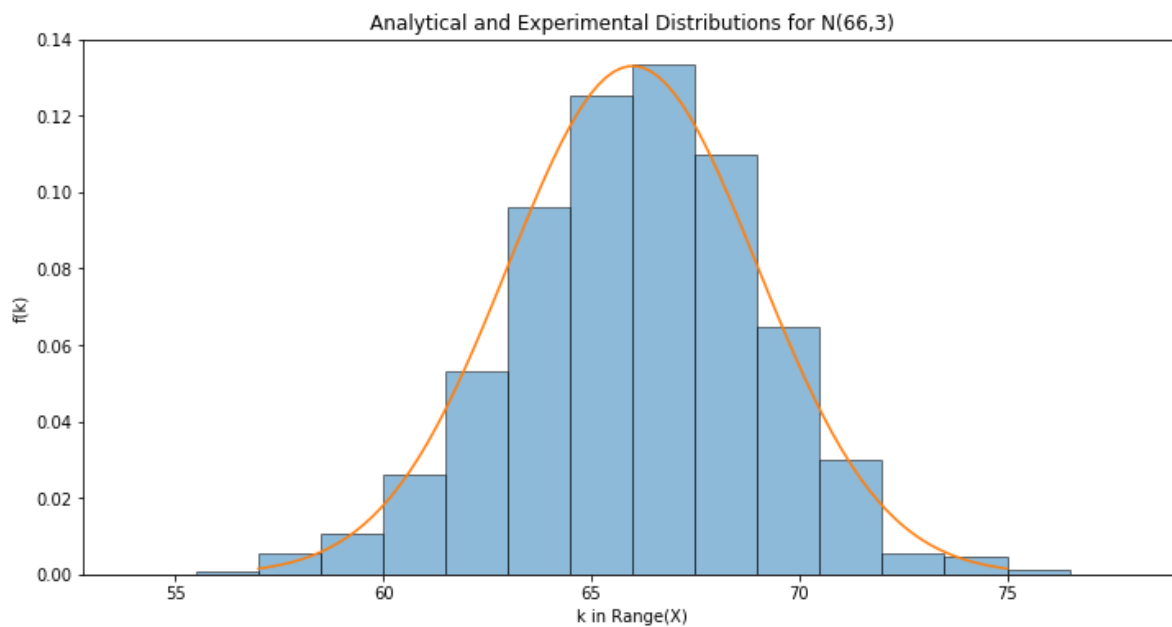
But now let's think about the issue of precision, i.e., the width of the bins. Again, try each of the following and see what happens. You can see that too-wide bins don't give much information, but too-narrow bins don't show how the data fits the normal distribution. There is a relationship between the number of data points and the width of the bins.

In [13]:

```
bin_width = 0.5          # try changing this to 0.5, 0.2, 0.1, 0.05, and 0.01

display_normal_samples_binned(66,3,1000,bin_width)
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:8: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.



Problem 11 (a)

Clearly the "fit" with the normal curve depends on the width of the bins! For the following three examples, find a value for the indicated parameter which gives a good correspondence between the normal curve and the data.

For the rest of these problems, just leave the code with your solution value and display the results -- no need to do anything else.

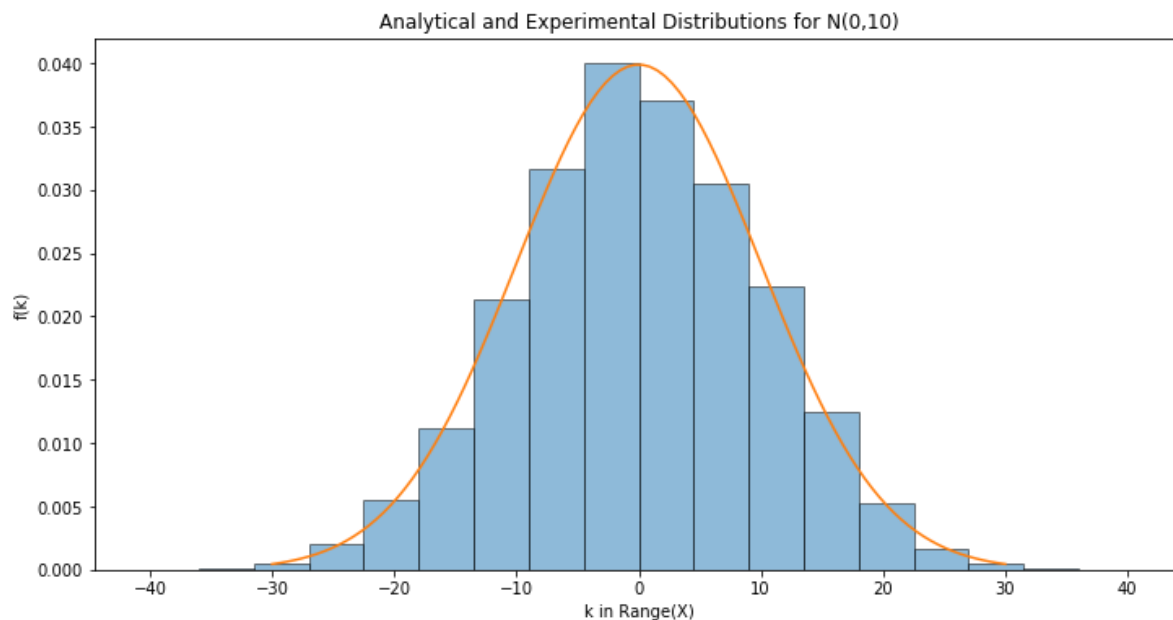
In [14]:

```
# Problem 1(a)

bin_width = 0.45    # experiment with this value 0.01, 0.05, 0.1, 0.15, etc. -- find
                    # largest number which still gives a good fit

display_normal_samples_binned(0,10,3000,bin_width)    # don't change this line
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:8: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.



Problem 11 (b)

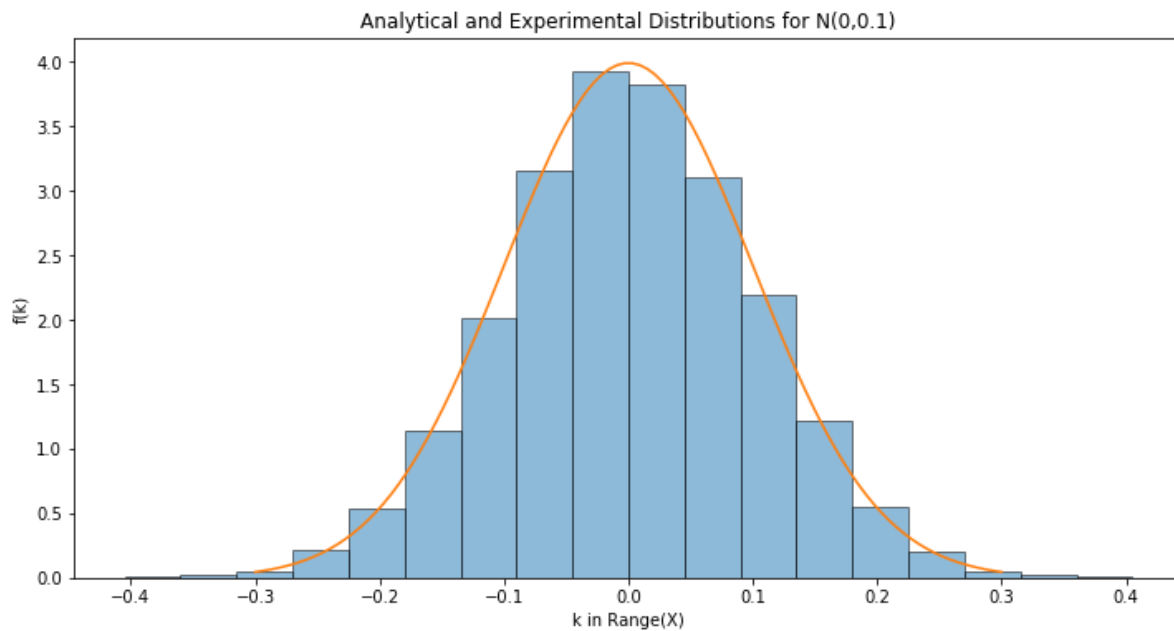
In [15]:

```
# Problem 1(b)

bin_width = 0.45      # experiment with this value -- find the largest number which s

display_normal_samples_binned(0,0.1,10000,bin_width)      # don't change this line
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:8: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.



Problem 11 (c)

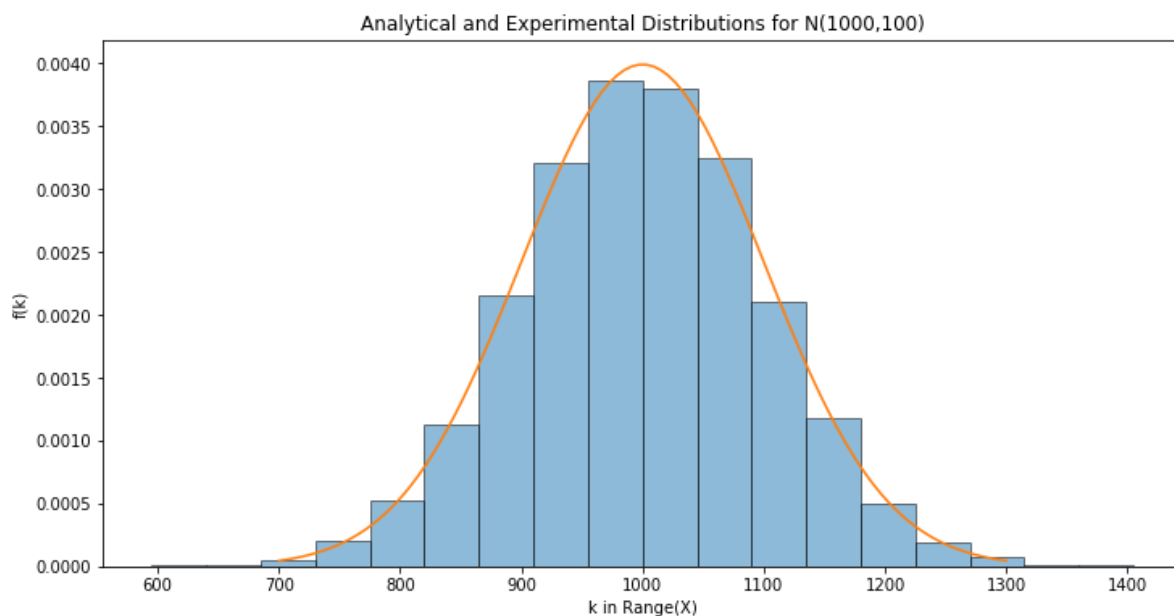
In [16]:

```
# Problem 11(c)

bin_width = 0.45      # experiment with this value -- find the largest number which st

display_normal_samples_binned(1000,100,10000,bin_width)      # don't change this line
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:8: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.



Problem 11 (d)

What do you think is a good value for `bin_width` in general, assuming that `num_trials` is sufficient to give a reasonable approximation of the normal distribution

Solution I believe that a good value for `bin_width` in general is 0.45.

Problem Twelve: Generating Poisson and Exponential Variates

Following up on HW 5, Problem 12, we now explore how to generate random variates from the Exponential and Poisson Distributions. We will display the results using a similar "binning" technique like we used in the previous problem.

Part (A)

The following formula should look familiar! If U is a random variable uniformly distributed in the range $[0..1)$, then

$$Y = \frac{-\ln(U)}{\lambda}$$

is a real number which is distributed according to the Exponential Distribution with rate parameter λ , i.e.,

$$Y \sim \text{Exp}(\lambda).$$

Note: \ln is log to the base e (just `log(...)` in Python).

For this problem, you must create a function `exp_rvs(...)` which generates exponential variates, using this formula (it is very similar to your solution to the last problem in HW 5), and demonstrate it by printing out 10 values for $\lambda = 4$.

In [17]:

```
## Solution for (A)

def exp_rvs(lam):
    value = random()
    y = -log(value)/lam
    return y

# test it

seed(0)

lam = 4

num_trials = 10

for x in range(num_trials):
    print(exp_rvs(lam))
```

```
0.04227577129370466
0.06928301239535861
0.21653514653972164
0.3378121740979418
0.16771205450973864
0.22600771216986748
0.060900798344498956
0.2982477266126604
0.18527102628047376
0.1347282520024191
```

Part (B)

Now you must complete the following code template to display the results of generating 10^5 exponential variates from $\text{Exp}(1)$, and finding a way to bin the results so that you get an accurate idea of how well they fit the theoretical distribution.

The binning here is a little simpler than in the case of the normal distribution, we will simply give the values we want to display (since the range is infinite) and the number of bins to show. This might involve a little tweaking to get a good diagram....

Optional: Can you think of a way of limiting the range based on λ ?

In [18]:

```

## Solution for (B)

def exp_pdf(lam,x):
    value = lam * 2.718281828 **(-1 * lam * x)
    return value

def display_exponential_samples_binned(lam,num_trials,limit,num_bins):
    fig, ax = plt.subplots(1,1,figsize=(12,6))
    plt.title('Analytical and Experimental Distributions for Exp('+str(lam)+'')
    plt.ylabel("f(k)")
    plt.xlabel("k in Range(X)")
    # use exp_rvs(...) to generate random samples
    X = [exp_rvs(lam) for k in range(num_trials)]
    # filter out those above the limit
    X1 = [ x for x in X if x < limit ]
    # make the bins
    bs = linspace(0,limit,num_bins)
    plt.hist(X,bins=bs,normed=True,edgecolor='k',alpha=0.5)
    #plt.hist(X,bins=bs,density=True,edgecolor='k',alpha=0.5) # use this 1.
    # Now generate the theoretical distribution with the same rate
    X2 = linspace(0,limit,100)
    Y = [ exp_pdf(lam,x) for x in X2]
    plt.plot(X2,Y)
    plt.show()

# test it

lam = 1

num_trials = 10**5

# You will only display values in range [0 .. limit]

limit = 1/lam + (1/lam) * 3 # Your answer here

# The range will be put into this many bins

num_bins = 50 # Your answer here

seed(0)

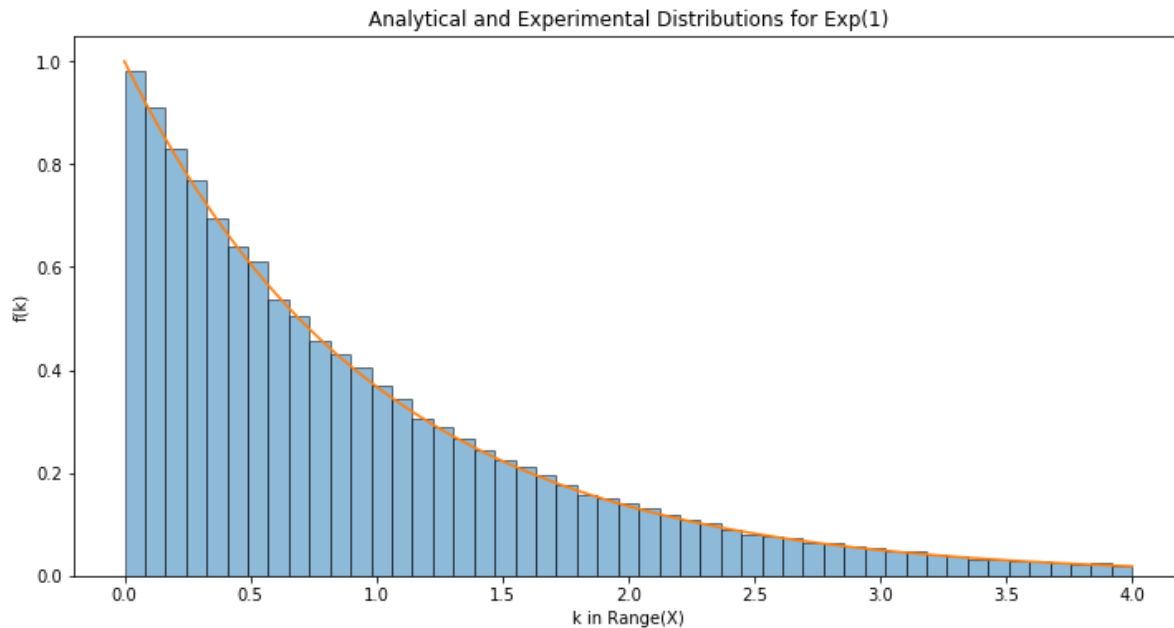
display_exponential_samples_binned(lam,num_trials,limit,num_bins)

```

```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:19: Matp
lotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be remove
d in 3.1. Use 'density' instead.

```



Part (C)

Now we will generate Poisson random variates by simulation, counting how many exponentially-distributed arrivals occur in a unit time (using our result from (a)).

For this problem, first recall that λ is the number of arrivals per unit time (meaning in 1.0 time units); so to simulate one "poke" at the Poisson RV `poi_rvs(lam)` you should "poke" the exponential RV you wrote above, summing the values until you exceed 1.0, then you can count how many arrivals in fact arrived between times 0 and 1.0.

Fill in the function stub and verify by running the code block to see the results (you can compare them with the theoretical distribution on the Distributions notebook).

In [19]:

```
## Solution for (C)
```

```
def poi_rvs(lam):
```

```
    sum_of = 0
```

```
    number = 0
```

```
    while sum_of <= 1:
```

```
        sum_of = sum_of + exp_rvs(lam)
```

```
        number = number + 1
```

```
    return number
```

```
# test it
```

```
lam = 4
```

```
num_trials = 10**6
```

```
seed(0)
```

```
Outcomes12c = [ poi_rvs(lam) for k in range(num_trials)]
```

```
show_distribution( Outcomes12c , "Empirical Probability Distribution for Poi(" + str(lam) + ")")
```

