

Statistical Methods in AI: Assignment1

Soumya Jahagirdar

September 2021

1 Introduction

This document contains the answers for the questions present in Assignment 1.

2 Questions

2.1 Question 1

Give an example each of probability mass functions with finite and infinite ranges. Show that the conditions on PMF are satisfied by your example.

Solution

Probability mass function of a random variable is the measure of probabilities of the possible set of values for a random variable, i.e

$$P_X(x_k) = P(X = x_k), k = 1, 2, 3, \dots, \quad (1)$$

An example of **Probability mass function** with **finite range** is an experiment of drawing a spade from a deck of cards. In this case, the task is to pick a card from a deck of cards and find the probability that a randomly picked card is a spade is

$$P(x = spade) = N_{spade}/N_{total} \quad (2)$$

where $P(x = spade)$ is the random experiment to pick a spade, and N_{spade} is the total number of spade cards in the deck and, N_{total} is the total number of cards present in the deck.

Similarly, along with $P(x = spade)$, if we consider $P(x = diamond)$, $P(x = club)$, $P(x = heart)$, then in general we can write the above equation as

$$P(x = v_i) = N_{v_i}/N_{total} \quad (3)$$

where v_i belongs to X and $X = P(x = diamond)$, $P(x = club)$, $P(x = heart)$, $P(x = spade)$.

Here we call $\mathbf{P}(\cdot)$ as the probability mass function of the experiment where a card is drawn from a deck of cards with finite range, i.e x takes a set of finite countable values.

In the above case probability of each case considered will be 0.25 as there are 13 cards of each type amongst a set of 52 cards. (assuming that there is replacement procedure followed).

An example of **Probability mass function** with **infinite range** is the experiment of obtaining a number divisible by 3 from a set of positive integers, i.e

Where $X = 3, 6, 9, 12, \dots$,

2.2 Question 2

Show with complete steps that the variance of uniform density is given by equation 10. (Hint: use the expression for variance in equation 5.)

Solution

Question 2 Solution

Solution

Variance of a density function is given by

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx, \text{ where } \mu \text{ is the mean and } p(x) \text{ is the probability of } x$$

We know that (from eqn (9)) and the value of μ for uniform density this function that,

$$p(x) = \frac{1}{(b-a)} \quad \& \quad \mu = \frac{(b+a)}{2}$$

\therefore The variance of Uniform density function will be as follows

$$\begin{aligned} \sigma^2 &= \int_a^b \left(x - \frac{(b+a)}{2} \right)^2 \left(\frac{1}{(b-a)} \right) dx \\ &= \int_a^b \left[x^2 - x(b+a) + \frac{(b+a)^2}{4} \right] \left[\frac{1}{(b-a)} \right] dx \\ &= \int_a^b \left[\frac{x^2}{(b-a)} - \frac{x(b+a)}{(b-a)} + \frac{(b+a)^2}{4(b-a)} \right] dx \\ &= \int_a^b \frac{x^2}{(b-a)} dx - \int_a^b \frac{x(b+a)}{(b-a)} dx + \int_a^b \frac{(b+a)^2}{4(b-a)} dx \\ &= \frac{1}{(b-a)} \left[\frac{x^3}{3} \right]_a^b - \frac{(b+a)}{(b-a)} \left[\frac{x^2}{2} \right]_a^b + \frac{1}{4} \frac{(b+a)^2}{(b-a)} [x]_a^b \\ &= \frac{1}{(b-a)} \left[\frac{b^3}{3} - \frac{a^3}{3} \right] - \frac{(b+a)}{(b-a)} \left[\frac{b^2}{2} - \frac{a^2}{2} \right] + \frac{1}{4} \frac{(b+a)^2}{(b-a)} (b-a) \\ &= \frac{(b-a)}{(b-a)} \frac{(a^2+ab+b^2)}{3} - \frac{(b+a)}{(b-a)} \frac{(b+a)(b-a)}{2} + \frac{1}{4} (a+b)^2 \\ &= \frac{4a^2+4ab+4b^2-6a^2-6b^2-12ab+3a^2+3b^2+6ab}{12} \\ &= \frac{a^2-2ab+b^2}{12} = \frac{(b-a)^2}{12} \longrightarrow (10) \end{aligned}$$

Figure 1: Solution for second question

2.3 Question 3

Show examples of two density functions (draw the function plots) that have the same mean and variance, but clearly different distributions. Plot both functions in the same graph with different colours.

Solution The following are the examples of the two density functions

1. A single normal distribution $N(0, 5)$
2. A combination of two normal distributions, where each distribution has standard deviation-1, and respective means are 2 and -2. The resulting combination has a mean of 0 and standard variance of 5.

```
In [5]: from scipy.optimize import curve_fit
from matplotlib import pyplot as plt
import numpy as np
x = np.linspace(-5,5,100)
mu = 0
std = 2.23
single = np.random.normal(mu, std)
mixture1 = np.random.normal(2, 1)
mixture2 = np.random.normal(-2, 1)
def funcfit(x1,x2,mu,sigma1,sigma2):
    return np.exp(-(x1-mu)**2/(2*sigma1**2))np.exp(-(x2-mu)**2/(2*sigma2**2))

init_guess=funcfit(mixture1,mixture2,2,-2,1,1)
count1, bins, count1 = np.histogram(x, bins=100)
pdf1 = count1 / sum(count1)

plt.plot(x,pdf1, color="red")
plt.plot(x,pdf2, color="green")
plt.show()
```

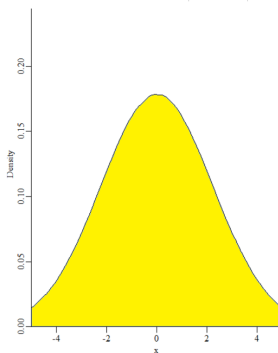


Figure 2: Solution for third question: example of one normal distribution function

These are clearly different distributions but they have same mean and variance.

```
In [5]: from scipy.optimize import curve_fit #non linear curve fitting tool
from matplotlib import pyplot as plt
import numpy as np
x = np.linspace(-1,5,10)
mu = 0
std = 2.23
single = np.random.normal(mu, std)
mixture1 = np.random.normal(2, 1)
mixture2 = np.random.normal(-2, 1)
def func2fit(x1,x2,m_1,m_2,std_1,std_2): #define a single gauss curve
    return np.exp(-(x1-m_1)**2/(std_1**2))+np.exp(-(x2-m_2)**2/(std_2**2))

init_guess=func2fit(mixture1,mixture2,2,-2,1,1)

count1, bins, counts = np.histogram(init_guess)
pdf1 = counts / sum(counts)

plt.plot(x,pdf1, color="red")
plt.plot(x,pdf2, color="green")
plt.show()
```

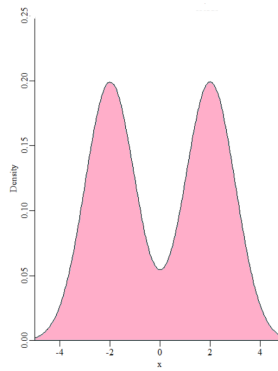


Figure 3: Solution for third question: example of two normal distribution function

2.4 Question 4

Show that the alternate expression for variance given in equation 5 holds for discrete random variables as well.

Solution

Solution 4 :

The alternate expression for variance given in equation (5) is

$$\sigma^2 = E[x^2] - (E[x])^2 \quad \text{---(5)}$$

Before this equation we know that

$$\sigma^2 = \int x^2 p(x) dx - 2\mu \int x p(x) dx + \mu^2 \int p(x) dx$$

As this equation was for continuous random variable, in case of discrete RV, the integrations get converted into summations i.e.,

$$\text{var}(x) = E[(x-\mu)^2] = \sum_{i=1}^n (x_i - \mu)^2 P(x_i)$$

when x is discrete RV

if w.r.t for x when discrete RV

$$E[x] = \sum_{x \in X} x p(x) \quad \text{or} \quad \sum_{i=1}^n x_i p(x_i)$$

$$E[x^2] = \sum_{x \in X} x^2 p(x) \quad \text{or} \quad \sum_{i=1}^n x_i^2 p(x_i)$$

$$\sigma^2 = E[(x-\mu)^2]$$

$$= \sum_{i=1}^n (x_i - \mu)^2 p(x_i)$$

$$= \sum_{i=1}^n (x_i^2 + \mu^2 + 2x_i\mu) p(x_i)$$

$$= \sum_{i=1}^n x_i^2 p(x_i) - 2\mu \sum_{i=1}^n x_i p(x_i) + \mu^2 \sum_{i=1}^n p(x_i)$$

$$= E[x^2] - (E[x])^2$$

Hence Proved //

Figure 4: Solution for fourth question

2.5 Question 5

Prove that the mean and variance of a normal density, $N(\mu, \sigma^2)$ are indeed its parameters, μ and σ^2 .

Solution

Question 5 Solution

→ Equation of normal density function

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

To find mean:

$$\mu = E(x) = \int_{-\infty}^{\infty} x f(x) dx = \int_{-\infty}^{\infty} x \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx$$

$$= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx$$

$\left\{ \begin{array}{l} \text{Let } \frac{x-\mu}{\sigma} = t \\ \text{or } dx = \sigma dt \end{array} \right\}$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\mu + \sigma t) e^{-\frac{t^2}{2}} \sigma dt$$

$$= \frac{1}{\sqrt{2\pi}} \left[\int_{-\infty}^{\infty} \mu e^{-\frac{t^2}{2}} dt + \int_{-\infty}^{\infty} \sigma t e^{-\frac{t^2}{2}} dt \right]$$

$\left[\begin{array}{l} \text{[even; even]} \\ \text{[odd; odd]} \end{array} \right]$

$$= \frac{2}{\sqrt{2\pi}} \int_0^{\infty} e^{-\frac{t^2}{2}} dt$$

$\left\{ \begin{array}{l} \text{Let } t^2/2 = p \\ \frac{2t}{2} \cdot dt = dp \Rightarrow dt = \frac{dp}{t} = \frac{dp}{\sqrt{2p}} \end{array} \right\}$

$$\therefore E(x) = \frac{2}{\sqrt{2\pi}} \mu \int_0^{\infty} e^{-p} \frac{dp}{\sqrt{2p}}$$

$$= \frac{2}{\sqrt{2\pi}} \mu \int_0^{\infty} p^{-1/2} e^{-p} dp$$

(integrating of $\int_0^{\infty} x^{n-1} e^{-x} dx = \text{gamma}(n)$)

$$= \frac{\mu}{\sqrt{\pi}} \text{gamma}(1/2)$$

$$= \frac{\mu}{\sqrt{\pi}} \times \sqrt{\pi}$$

$$E(x) = \mu$$

Hence Proved that 'mean' of the normal density is its parameter ' μ '.

Figure 5: Solution for fifth question "mean"

Questions Solution for Variance

For variance we need $E(x^2)$

$$E(x^2) = \int_{-\infty}^{\infty} x^2 f(x) \cdot dx$$

$$= \int_{-\infty}^{\infty} x^2 \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \cdot dx$$

$$= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \cdot dx$$

Let $\frac{x-\mu}{\sigma} = t \Rightarrow x = \mu + \sigma t$
 $dx = \sigma dt$

$$= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} (\mu + \sigma t)^2 e^{-\frac{1}{2}t^2} \cdot \sigma dt$$

$$= \frac{1}{\sqrt{2\pi}} \left[\int_{-\infty}^{\infty} \mu^2 e^{-\frac{1}{2}t^2} \cdot \sigma dt + \int_{-\infty}^{\infty} 2\mu\sigma t e^{-\frac{1}{2}t^2} \cdot \sigma dt + \int_{-\infty}^{\infty} \sigma^3 t^2 e^{-\frac{1}{2}t^2} \cdot dt \right]$$

Let $\frac{t^2}{2} = p$
 $\frac{2t}{2} \cdot dt = dp \Rightarrow dt = \frac{dp}{t} = \frac{dp}{\sqrt{2p}}$

$$E(x^2) = \frac{2\mu^2}{\sqrt{2\pi}} \int_0^{\infty} e^{-p} \cdot \frac{1}{\sqrt{2p}} \cdot \frac{dp}{\sqrt{2p}} + \frac{2\sigma^2}{\sqrt{2\pi}} \int_0^{\infty} p e^{-p} \cdot \frac{1}{\sqrt{2p}} \cdot \frac{dp}{\sqrt{2p}}$$

$$= \frac{2\mu^2}{\sqrt{2\pi}} \int_0^{\infty} e^{-p} \frac{dp}{2p} + \frac{2\sigma^2}{\sqrt{2\pi}} \int_0^{\infty} p e^{-p} \frac{dp}{2p}$$

$$= \frac{\mu^2}{\sqrt{\pi}} \int_0^{\infty} p^{-1/2} e^{-p} \cdot dp + \frac{\sigma^2}{\sqrt{\pi}} \int_0^{\infty} p^{1/2} e^{-p} \cdot dp$$

$$= \frac{\mu^2}{\sqrt{\pi}} \text{gamma}(1/2) + \frac{2\sigma^2}{\sqrt{\pi}} \text{gamma}(3/2)$$

$$= \frac{\mu^2}{\sqrt{\pi}} \times \frac{\sqrt{\pi}}{2} + \frac{2\sigma^2}{\sqrt{\pi}} \times \frac{\sqrt{\pi}}{2}$$

$$E(x^2) = \mu^2 + \sigma^2$$

$$\text{var}(x) = [E(x^2)] - [E(x)]^2$$

$$= \mu^2 + \sigma^2 - \mu^2 = \sigma^2$$

Hence proved that "standard deviation" of the normal density is its parameter (σ^2).

Figure 6: Solution for fifth question "variance"

2.6 Question 6

Using the inverse of CDFs, map a set of 10,000 random numbers from $U[0; 1]$ to follow the following pdfs:

1. Normal density with $\mu = 0$, $\sigma = 3.0$.
2. Rayleigh density with $\sigma = 1.0$.
3. Exponential density with $\lambda = 1.5$

Once the numbers are generated, plot the normalized histograms (the values in the bins should add up to 1) of the new random numbers with appropriate bin sizes in each case; along with their pdfs. What do you infer from the plots? Note: see `rand()` function in C for $U[0, \text{INT_MAX}]$. **Solution**

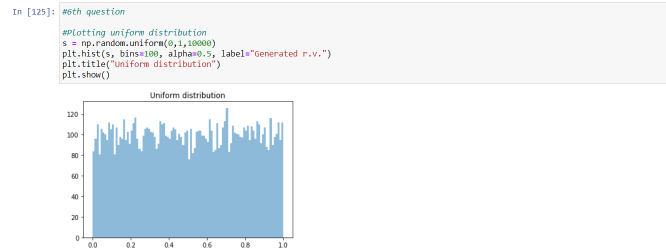


Figure 7: Solution for sixth question uniform distribution

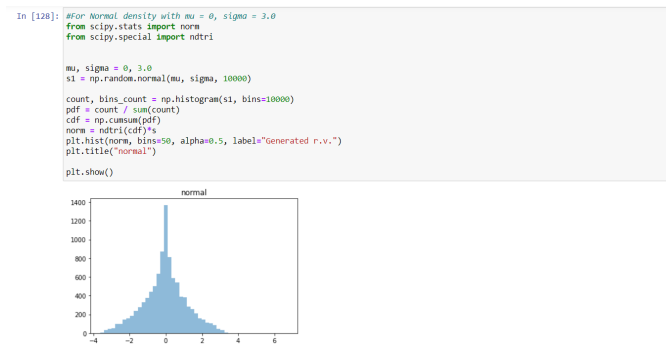


Figure 8: Solution for sixth question normal distribution

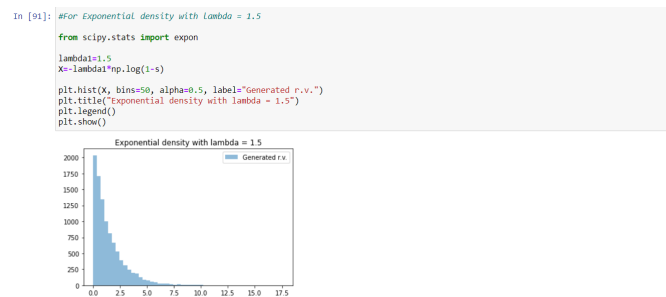


Figure 9: Solution for sixth question exponential distribution

2.7 Question 7

Write a function to generate a random number as follows: Every time the function is called, it generates 500 new random numbers from $U[0, 1]$ and outputs their sum.

Generate 50,000 random numbers by repeatedly calling the above function, and plot their normalized histogram (with bin-size = 1). What do you find about the shape of the resulting histogram?

Solution

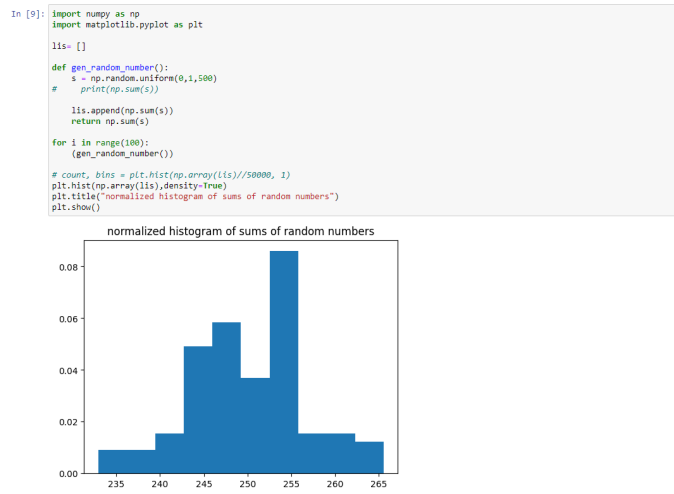


Figure 10: Solution for seventh question

After adding the numbers sampled from randomly distribution, the resultant of the sum of the numbers obtained is not uniformly distribution. It is seen that the numbers are normally distributed.