

# Advanced Machine Learning: Exercises

Ramon Ruiz Dolz

April 2019

## Contents

- 1 **Exercise(\*):** *Compute  $H(C)$ ,  $H(S)$ ,  $H(C|S)$ ,  $H(S|C)$ ,  $H(C, S)$ , and  $I(C; S)$ .* 1
  - 2 **Exercise(\*\*):** *Reproduce an example similar to the previous example with three unidimensional distributions, with known mean and equal and known variance where only  $\pi_1$ ,  $\pi_2$ ,  $\pi_3$  are unknown. The data have to be generated with three gaussians. Obtain the corresponding plots.* 2
  - 3 **Exercise(\*\*\*):** *Reproduce an example similar to the previous example with two uni-dimensional distributions, with equal and known variance where the mean of each distribution and  $\pi_1$ ,  $\pi_2$  are unknown. Obtain the corresponding plots.* 4
  - 4 **Exercise(\*\*\*):** *Reproduce an example similar to the previous example with two uni-dimensional distributions, with equal and known mean where the variances and  $\pi_1$ ,  $\pi_2$  are unknown. Obtain the corresponding plots.* 6
- 
- 1 **Exercise(\*):** *Compute  $H(C)$ ,  $H(S)$ ,  $H(C|S)$ ,  $H(S|C)$ ,  $H(C, S)$ , and  $I(C; S)$ .*

En este ejercicio se nos pide realizar los cálculos de la entropía de las variables Climatología (C) y Luminosidad (L) así como las entropías condicionales, la entropía conjunta y la intersección para estas mismas variables. Como punto de partida se toman los datos presentados en la diapositiva 13. Para la variable C, sus estados pueden ser, *despejado* (DES), *nublado* (NUB) y *lluvia* (LLU). Por otra parte, para la variable S, los estados posibles son desplazamiento *seguro* (SEG) o con *accidente* (ACC). A partir de las probabilidades conjuntas proporcionadas, es posible obtener las probabilidades de cada variable marginalizando. Por lo tanto obtenemos las siguientes probabilidades a priori para la variable C:

- $P(\text{DES}) = 0.46$
- $P(\text{NUB}) = 0.33$
- $P(\text{LLU}) = 0.21$

Y para la variable S:

- $P(\text{SEG}) = 0.86$
- $P(\text{ACC}) = 0.14$

Una vez obtenidas las distintas probabilidades a priori, ya es posible realizar los cálculos requeridos para el ejercicio. Mediante la formula de la entropía:

$$H(X) = - \sum_{i=1}^{|X|} p(i) \log(p(i)) \quad (1)$$

Por lo tanto, se obtiene que:

$$H(C) = -(0.46\log(0.46) + 0.33\log(0.33) + 0.21\log(0.21)) \simeq 0.468 \quad (2)$$

$$H(S) = -(0.86\log(0.86) + 0.14\log(0.14)) \simeq 0.176 \quad (3)$$

A demás de las entropías de cada variable, también se nos pide calcular las entropías condicionales de estas dos variables. Para ello hacemos uso de la siguiente fórmula:

$$H(Y|X) = - \sum_x \sum_y p(x,y) \log(p(y|x)) \quad (4)$$

Que aplicada a nuestras variables nos conduce a los siguientes cálculos:

$$\begin{aligned} H(C|S) = & \quad (5) \\ & -((0.43\log(0.5) + 0.3\log(0.35) + 0.13\log(0.15) + (0.03\log(0.21) + 0.03\log(0.21) + 0.08\log(0.57))) = \\ & -((-0.373) + (-0.042)) \simeq 0.415 \end{aligned}$$

$$\begin{aligned} H(S|C) = & -((0.43\log(0.93) + 0.03\log(0.07)) + (0.3\log(0.91) + 0.03\log(0.09)) \\ & + (0.13\log(0.62) + 0.08\log(0.38))) = -((-0.048) + (-0.044) + (-0.061)) \simeq 0.153 \end{aligned} \quad (6)$$

Una vez obtenidas tanto las entropías como las entropías condicionales asociadas a cada variable, es posible finalizar los cálculos de la entropía conjunta y la interacción. Para ello se haran uso de las siguientes dos ecuaciones:

$$H(X, Y) = H(X) + H(Y|X) \quad (7)$$

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \end{aligned} \quad (8)$$

Con lo que se obtiene:

$$\begin{aligned} H(C, S) &\simeq 0.468 + 0.153 \\ &= 0.621 \end{aligned} \quad (9)$$

$$\begin{aligned} I(C; S) &\simeq 0.468 - 0.415 \\ &= 0.053 \end{aligned} \quad (10)$$

**2 Exercise(\*\*):** *Reproduce an example similar to the previous example with three unidimensional distributions, with known mean and equal and known variance where only  $\pi_1$ ,  $\pi_2$ ,  $\pi_3$  are unknown. The data have to be generated with three gaussians. Obtain the corresponding plots.*

En este ejercicio se nos pide reproducir el experimento mostrado en clase relacionado con el algoritmo EM. Para ello se han generado 500 muestras unidimensionales de forma aleatoria siguiendo tres distribuciones normales distintas con la varianza predefinida (2) y las medias (0, 1 y 2). Las muestras generadas se pueden observar en la Figure 1.

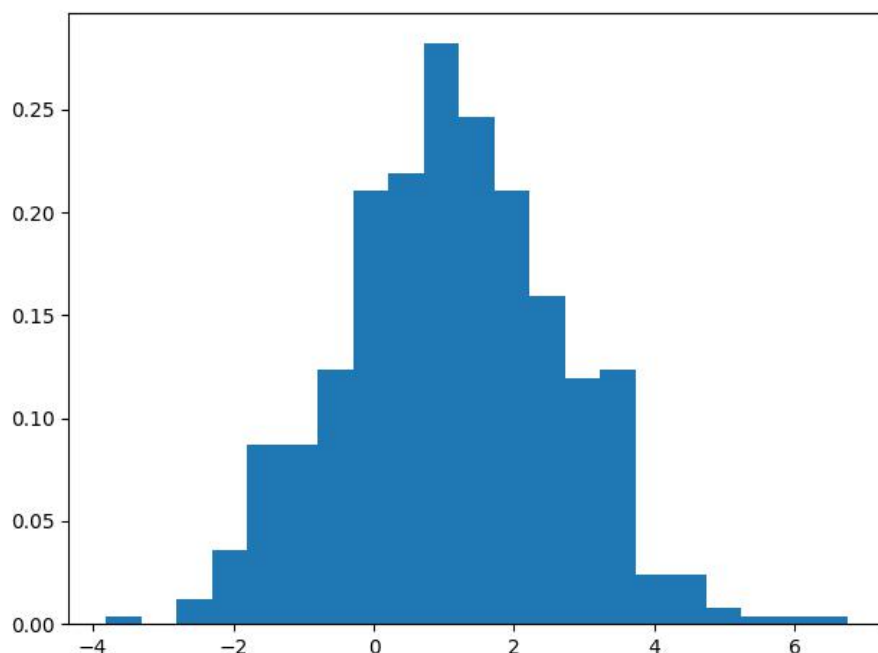


Figure 1: Distribución de las 500 muestras generadas aleatoriamente

Sobre este conjunto de muestras se ha aplicado el algoritmo EM con el objetivo de hallar los tres parámetros  $\pi$  que definen la mixtura de Gaussianas. Al finalizar este algoritmo se obtiene la estimación de los parámetros  $\pi_1$  (0.185),  $\pi_2$  (0.510) y  $\pi_3$  (0.305). En la Figure 2 se puede observar gráficamente la mixtura sobre la distribución de las muestras generadas.

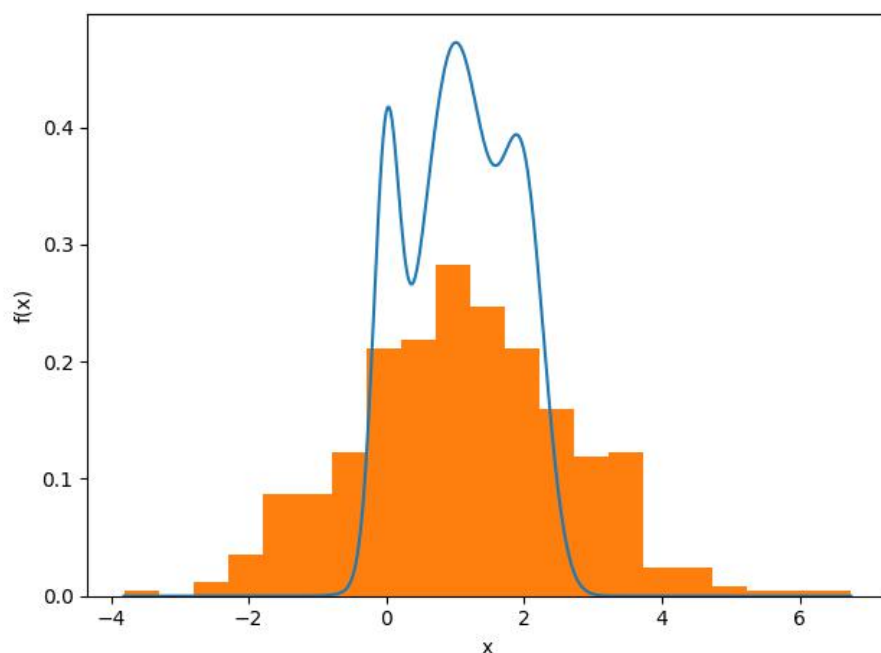


Figure 2: Mixtura de Gaussianas aprendida a partir de las muestras

Además de este primer experimento se ha lanzado un segundo experimento con 100 muestras con tal de aportar riqueza al ejercicio. En este segundo experimento se han generado 100 muestras unidimensionales de forma aleatoria siguiendo tres distribuciones normales distintas con la varianza predefinida (2) y las medias (0, 1 y 2). Las muestras generadas se pueden observar en la Figure 3.

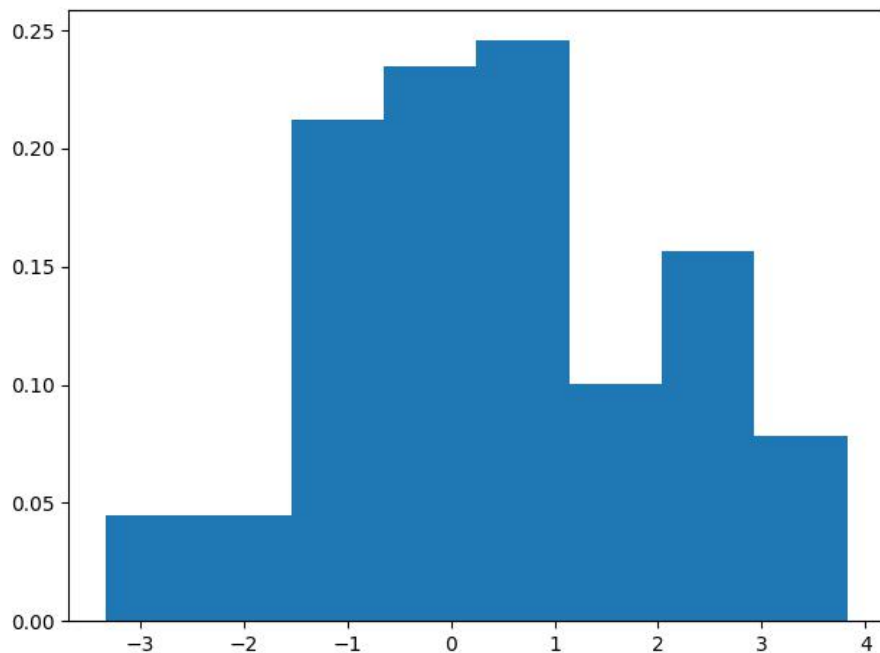


Figure 3: Distribución de las 100 muestras generadas aleatoriamente

Sobre este conjunto de muestras se ha aplicado el algoritmo EM con el objetivo de hallar los tres parámetros  $\pi$  que definen la mixtura de Gaussianas. Al finalizar este algoritmo se obtiene la estimación de los parámetros  $\pi_1$  (0.632),  $\pi_2$  (0.263) y  $\pi_3$  (0.104). En la Figure 4 se puede observar gráficamente la mixtura sobre la distribución de las muestras generadas.

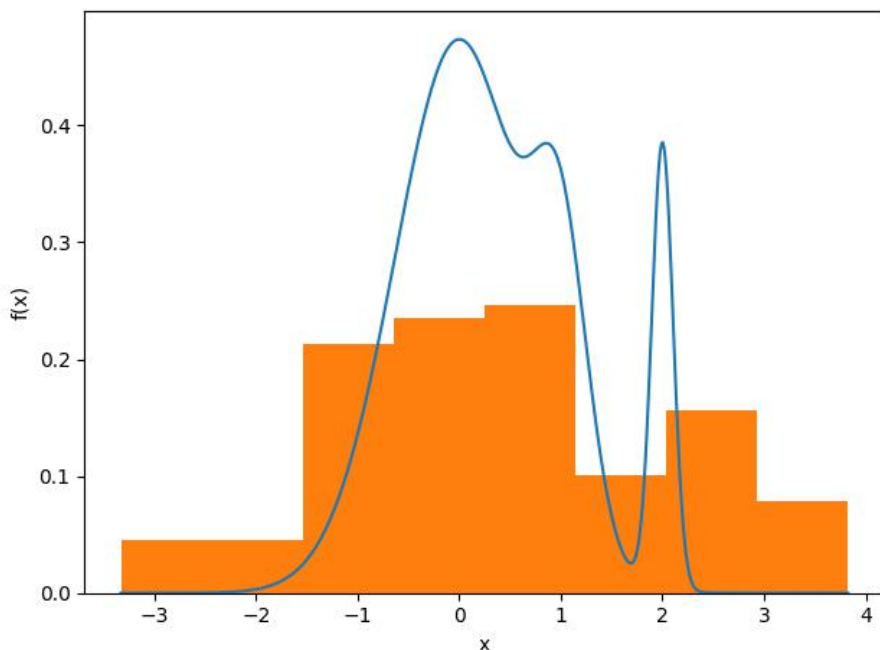


Figure 4: Mixtura de Gaussianas aprendida a partir de las muestras

### 3 Exercise(\*\*): *Reproduce an example similar to the previous example with two uni-dimensional distributions, with equal and known variance where the mean of each distribution and $\pi_1$ , $\pi_2$ are unknown. Obtain the corresponding plots.*

En este ejercicio se nos pide reproducir el experimento mostrado en clase relacionado con el algoritmo EM. Para ello se han generado 100 muestras unidimensionales de forma aleatoria siguiendo dos distribuciones normales distintas con la varianza predefinida (2). Las muestras generadas se pueden observar en la Figure 5.

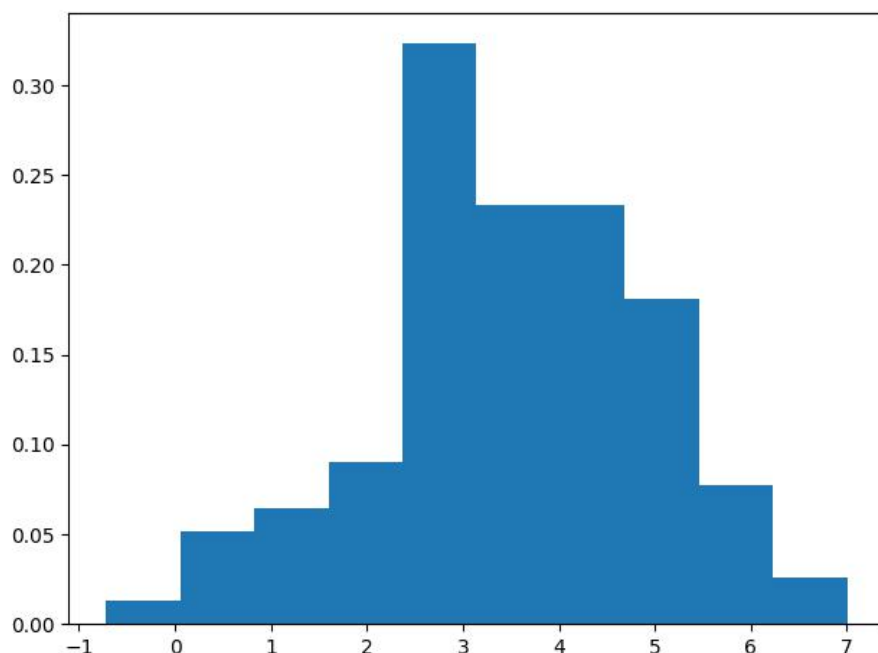


Figure 5: Distribución de las 100 muestras generadas aleatoriamente

Puesto que tanto los pesos como las medias son desconocidas, se inicializan aleatoriamente al empezar el algoritmo EM. Mediante la iteración de este algoritmo, a partir de los datos conocidos (distribución de las muestras y las varianzas) se trata de aprender la distribución de la mezcla de ambas Gaussianas. Al finalizar este algoritmo se estiman las medias (2.9 y 3.7) así como los parámetros  $\pi_1$  (0.67) y  $\pi_2$  (0.33). En la Figure 6 se puede observar gráficamente este resultado.

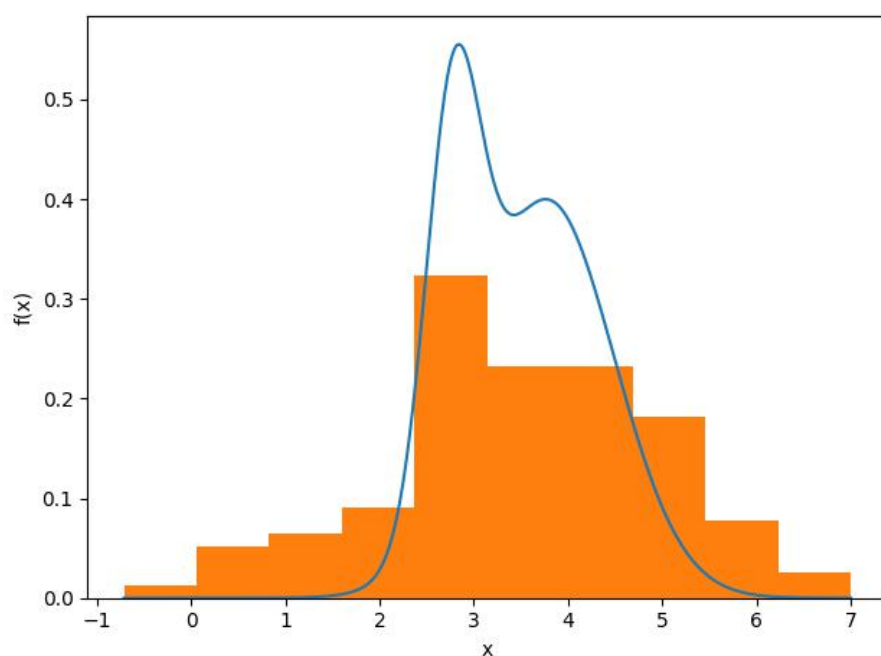


Figure 6: Mixtura de Gaussianas aprendida a partir de las muestras

Como bien se puede apreciar, los parámetros aprendidos que definen a cada Gaussiana no son iguales y, por lo tanto, la distribución de probabilidades tras aplicar el EM sigue una distribución similar al de la naturaleza de las muestras.

Debido a la aleatoriedad del ejercicio se ha realizado un segundo experimento con el objetivo de aportar mayor riqueza. En este segundo experimento se han vuelto a generar 100 muestras unidimensionales de forma aleatoria siguiendo dos distribuciones Gaussianas distintas con varianzas conocida a priori (2). En la Figure 7 se puede observar la distribución del conjunto total de muestras generadas.

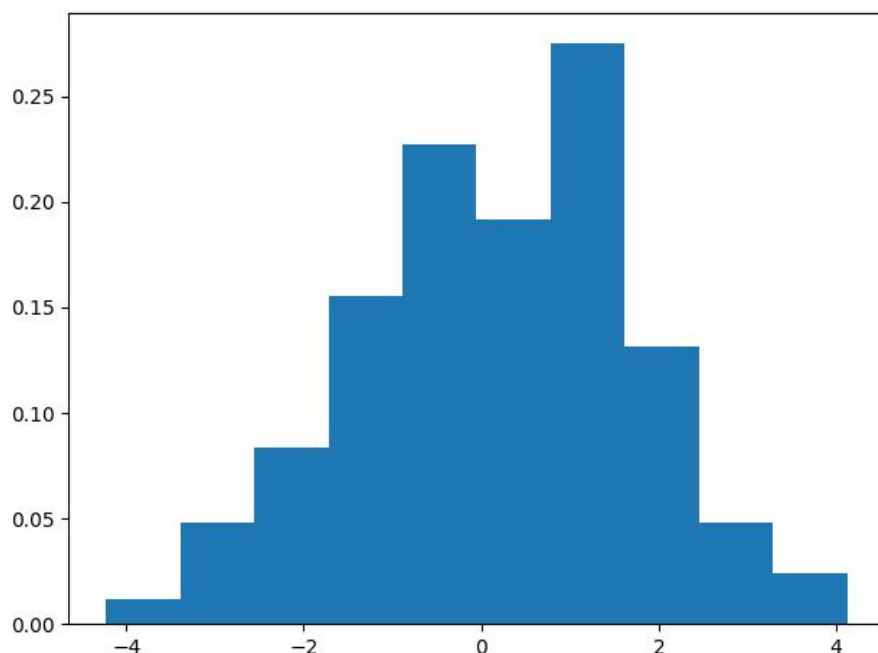


Figure 7: Distribución de las 100 muestras generadas aleatoriamente

Del mismo modo que en el ejercicio anterior, las medias y los pesos son desconocidos. Al finalizar el algoritmo, se estiman las medias (-0.5 y 0.8) así como los parámetros  $\pi_1$  (0.498) y  $\pi_2$  (0.502) obteniendo como resultado la mixtura observable en la Figure 8

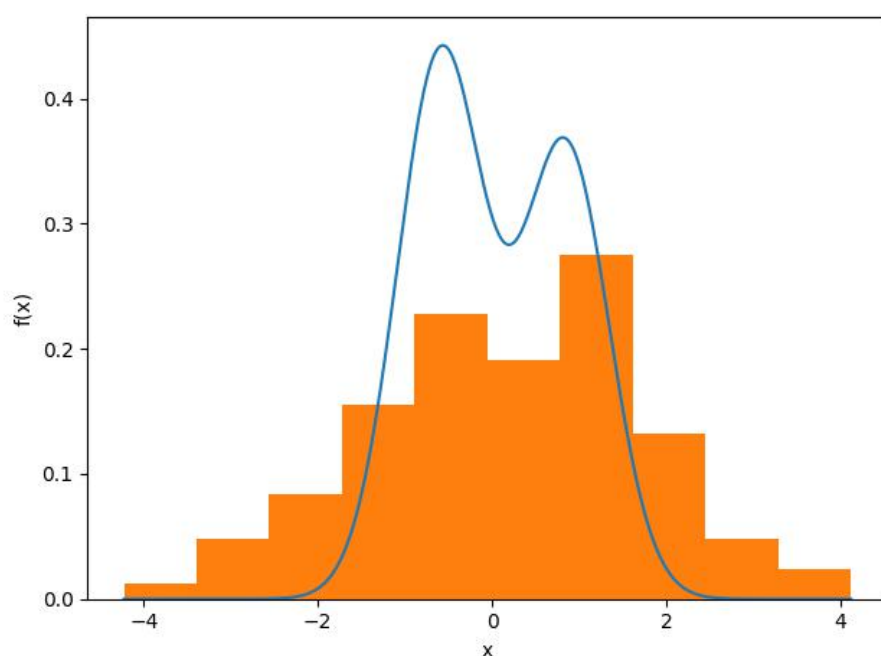


Figure 8: Mixtura de Gaussianas aprendida a partir de las muestras

**4 Exercise(\*\*): *Reproduce an example similar to the previous example with two uni-dimensional distributions, with equal and known mean where the variances and  $\pi_1$ ,  $\pi_2$  are unknown. Obtain the corresponding plots.***

En este ejercicio se nos pide reproducir el mismo experimento que en el anterior, pero en este caso las medias son conocidas y las varianzas desconocida. Para ello, de forma similar al ejercicio anterior se han realizado dos experimentos con datos generados aleatoriamente distintos.

En el primer experimento se han generado 100 muestras unidimensionales siguiendo dos distribuciones normales distintas con medias conocidas (0 y 1). Podemos observar la distribución de las muestras generadas en este experimento en Figure 9.

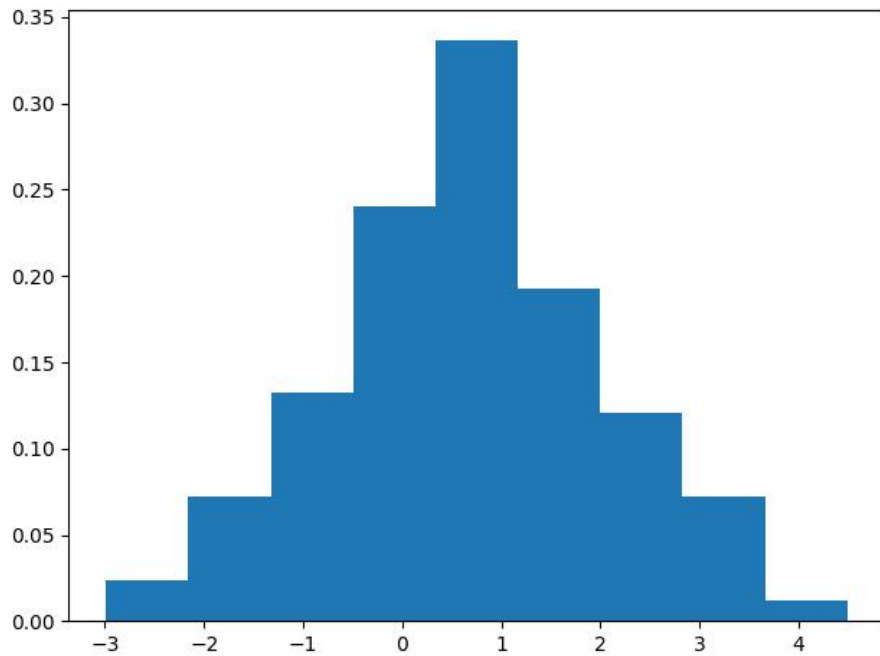


Figure 9: Distribución de las 100 muestras generadas aleatoriamente

En este caso, las varianzas y los pesos son desconocidos. Mediante el uso del algoritmo EM se han conseguido estimar las varianzas (1.47 y 1.91) y los parámetros  $\pi_1$  (0.61) y  $\pi_2$  (0.39). Las mixturas de Gaussianas resultantes se pueden observar gráficamente en la Figure 10

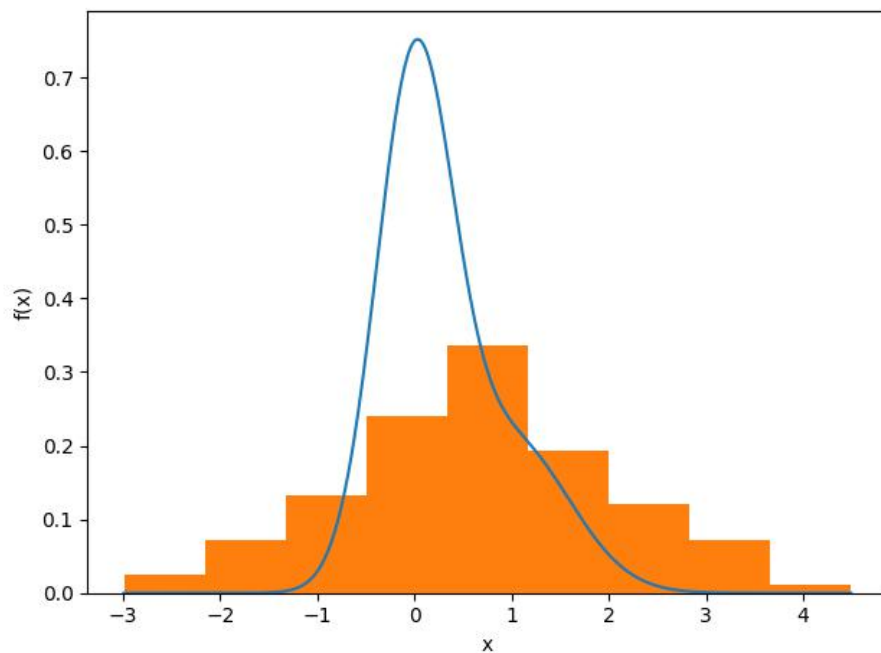


Figure 10: Mixtura de Gaussianas aprendida a partir de las muestras

Finalmente, con tal de aportar riqueza al ejercicio se ha añadido un segundo experimento. De manera muy similar al experimento explicado anteriormente, se han generado 100 muestras unidimensionales de forma aleatoria con medias conocidas (0 y 1), siguiendo dos distribuciones Gaussianas diferentes. En la Figure 11 se puede apreciar la distribución de las 100 muestras generadas.

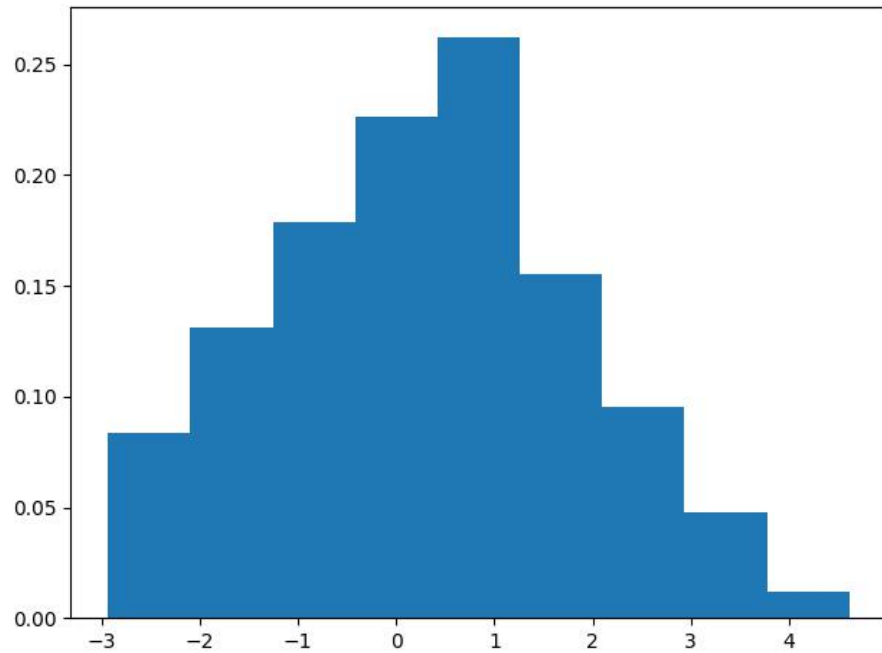


Figure 11: Distribución de las 100 muestras generadas aleatoriamente

A partir de estos datos generados se ha aprendido una mixtura de Gaussianas haciendo uso del algoritmo EM. Mediante este algoritmo se han estimado las varianzas (2.19 y 2.12) y los parámetros  $\pi_1$  (0.325) y  $\pi_2$  (0.675). El modelo aprendido se puede observar en la Figure 12.

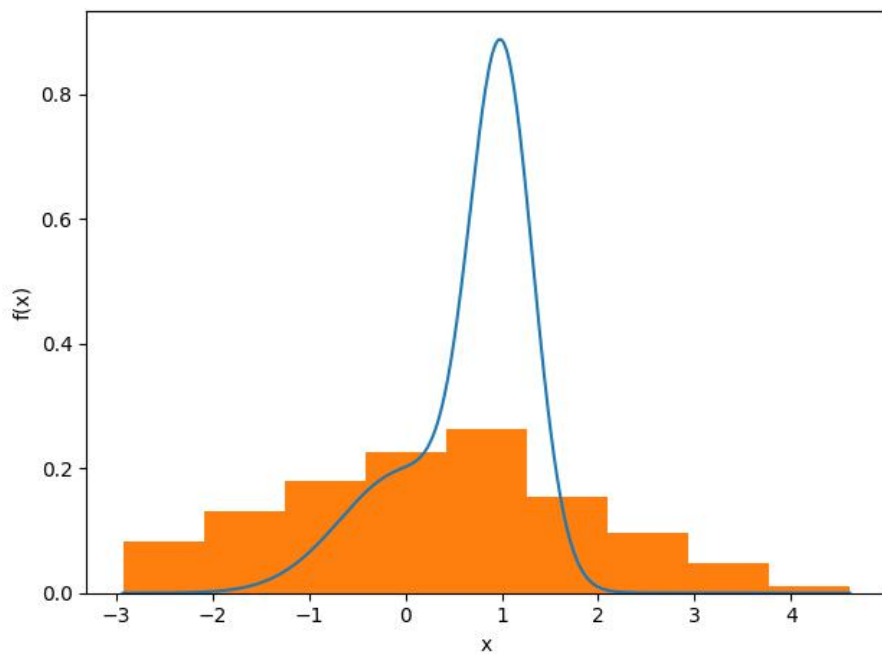


Figure 12: Mixtura de Gaussianas aprendida a partir de las muestras