

Proyecto Final de Programación Básica: Cadenas de Markov Chain Montecarlo con el método Metrópolis.

Universidad de Guanajuato, División de Ciencias e Ingenierías
Programación Básica
Profesora: Dra. Alma Xochitl González Morales
Alumnos: Guadalupe Sinaí Florián Landa, Ian Barrera Salas

December 8, 2018

1 Descripción del Programa

El método de ajuste de parámetros conocido como minimización de X^2 es un caso particular. En el programa y_i^{obs} y x_i^{obs} son un conjunto de coordenadas recolectadas de un experimento de péndulo oscilatorio donde x =el periodo de oscilación τ y y =el largo del cable l , los cuales tienen su error σ_x σ_y respectivamente. Para encontrar el valor de la pendiente ocupamos un ajuste de mínimos cuadrados expresada por la ecuación:

$$X^2 = \sum_i \frac{(y_i^{obs} - y_i^{mod}(x_i, m, b))^2}{\sigma_i^2} \quad (1)$$

El problema en el programa a resolver es la probabilidad de que cierto valor de la pendiente sea el valor verdadero que se ajuste a los datos. Esto se puede resolver con el teorema de Bayes:

$$P(\vec{\theta} | \vec{y}) \propto L(\vec{y} | \vec{\theta}) * Pr(\vec{\theta}) \quad (2)$$

$$P(A_i/B) = \frac{P(A_i) \cdot P(B/A_i)}{P(B)}$$

Donde:

$P(A_i)$ = Probabilidad a priori

$P(B/A_i)$ = Probabilidad condicional

$P(B)$ = Probabilidad Total

$P(A_i/B)$ = Probabilidad a posteriori

Figure 1: Ecuación del teorema de Bayes, especificando el significado de cada variable.

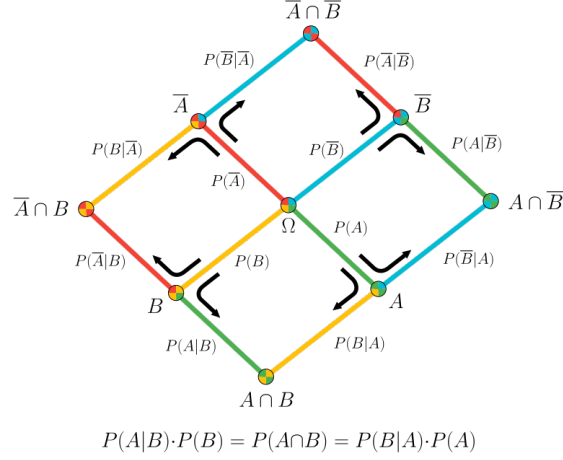


Figure 2: Visualización del treorema de Bayes por superposición de dos árboles de decisión.

P es la distribución de probabilidad prosterior, L es Likelihood o probabilidad de de verosimilitud, Pr es la distribución a priori. θ representa los parámetros a ajustar en el modelo inicial. Lo que se busca es encontrar el máximo de distribución posterior, es de decir, los parámetros que suponen mayor probabilidad de describir correctamente los datos iniciales $\vec{\theta} = (m, b)$ que son la pendiente y ordenada al origen. Lo más útil en este caso es suponer una distribución Gaussiana para L que se relacione con X^2 :

$$\ln(P(y|m)) = 0.5 \sum_i ((y_i^{obs} - y_i^{mod}(x_i, m)) - \sigma_i^2) \quad (3)$$

Por lo que el teorema de Bayes se expresaría como:

$$\ln(P(y|m)) \approx \ln(L(\vec{y}|\vec{\theta})) + \ln(Pr(\vec{y}|\vec{\theta})) \quad (4)$$

También se usará la distribución a priori plana, considerando los parámetros m,b y su probabilidad a priori de describir los datos en intervalos $m_{min} < m < m_{max}$ y $b_{min} < b < b_{max}$, si lo anterior se cumple el valor a priori será igual a zero. El método montercarlo de maximización es Metropolis-Hasting, que genera parámetros m,b de forma aleatoria y calcula la L a priori de dichos parámetros. Cada conjunto nuevo que se genera depende del valor de L y que tan mayor o menor al anterior sea.

A continuación se describirá el conjunto de programas escritos en lenguaje C, que nos permitió encontrar la pendiente con los datos (x,y) y error dados desde un archivo.

2 Comandos usados en el Programa

perce.h

Este programa nos ayudó a crear una librería con funciones que se definirán en el siguiente programa escrito:

Primero se indicaron las librerías que se usarán en el programa, las cuales incluyen definiciones comunmente usadas en algoritmos. En este programa se usaron: dos librerías estandar y la librería con varias funciones matemáticas:

$$\textit{include} < \textit{stdio.h} > \quad (5)$$
$$\textit{include} < \textit{math.h} > \quad (6)$$
$$\textit{include} < \textit{stdlib.h} > \quad (7)$$

float nos ayuda a declarar variables con valores de punto flotante, percentx() y percenty() nos ayudan a declarar las funciones que estarán en la librería perce.h

$$\textit{float percenty} (\textit{float} \textit{m0}, \textit{float} \textit{b0}); \quad (8)$$
$$\textit{float percentx} (\textit{float} \textit{m0}, \textit{float} \textit{b0}); \quad (9)$$

func.c

Este programa nos ayudó a definir las acciones que realizan las funciones que se encuentran dentro de la librería perce.h:

Se declara la función percentx(), la cual es sin argumentos de entrada pero con argumentos de salida . El nombre de la función debe ser único y exclusivo para dicha función.

$$\textit{float percentx} (\textit{float} \textit{m0}, \textit{float} \textit{b0}) \quad (10)$$

int, float y FILE ayudan a declarar las variables de tipo entero, punto flotante y archivo respectivamente. result [10] crea un arreglo con capacidad de 10 datos que se podrán guardar dentro de él. *xobs crea una variable de tipo apuntador con tamaño de float.

$$\textit{inti}, \textit{j}, \textit{n} = 10; \quad (11)$$
$$\textit{float} * \textit{xobs}, * \textit{yobs}, * \textit{sigmam}, * \textit{sigmab}, \textit{result}[10]; \quad (12)$$
$$\textit{FILE} * \textit{datos}; \quad (13)$$

Se reserva el espacio en la memoria usando malloc, que asigna un número de bytes indicados y devuelve un apuntador al primer byte del espacio asignado (un solo bloque), y se le indica que el total de bytes son el valor num por el tamaño de una variable entera

$$xobs = (float*)malloc(10 * sizeof(float)); \quad (14)$$

Abro el archivo dónde se encuentran los datos de (x,y):

$$info = fopen ("datos", "r"); \quad (15)$$

Se inicia un ciclo dónde el contador de las coordenadas i inicia en cero, la condición de seguir realizando el ciclo es que se repita un un numero menor al valor de la variable "n" y que al coontador se le aumente una unidad cada vez que se completa la serie de instrucciones. fopen lee los datos ya formateados del stdin y los guarda los valroes en las variables indicadas con "":

$$for(i = 0; i < n; i++) \quad (16)$$

$$fscanf(datos, "%f ", xobs[i]); \quad (17)$$

Cierro el archivo con datos:

$$fclose(datos); \quad (18)$$

Se inicia un ciclo dónde el contador de las coordenadas i inicia en cero, la condición de seguir realizando el ciclo es que se repita un un numero menor al valor de la variable "n" y que al coontador se le aumente una unidad cada vez que se completa la serie de instrucciones.

$$for (i = 0; i < 10; i++) \quad (19)$$

Se calcula la L=likelihood por el método metrópolis, el cual indica que el arreglo xobs[i] se multiplica por la m0 y se suma a b0, y luego se resta al arreglo yobs [i], esto se guarda en result [i]. En la segunda los arreglos de sigma se multiplican y se restan con la raíz del arreblo result [i], dicho valor es cuardado den result [i]. Se procede a relaizar una sumatoria con los valores del arreglo result y se multiplicar por -0.5. De dicho resultado se obtiene la exponencial que se guarda en res2.

$$result[i] = yobs[i] - (xobs[i] * (m0)) + b0; \quad (20)$$

$$result[i] = pow(result[i], 2) - (sigmam[i] * sigmam[i]); \quad (21)$$

$$suma+ = result[i]; \quad (22)$$

$$suma* = -0.5; \quad (23)$$

$$res2 = pow(2.718282, suma); \quad (24)$$

La segunda función es percenty, que se asemeja con con la función percentx, la única diferencia es que en la función percenty no se obtiene el exponencial de su respectiva sumatoria.

mainp.c

Este es el programa principal que calculará la pendiente m y la ordenada b que más se aproxime de acuerdo con las coordenadas (x,y) del experimento.

Se declaran las variables int, float y double que definen variables de tipo archivo, entero, punto flotante con 6 cifras y punto flotante con 12 cifras respectivamente.

$$FILE * resultados; \quad (25)$$

$$int i, j, n = 10; \quad (26)$$

$$float m0, b0, mn, bn; \quad (27)$$

$$double px0, pxn = 0, pxn2 = 0, px02, ale; \quad (28)$$

Se le pide al usuario, imprimiendo en la pantalla con printf, entrar los valores de la pendiente y ordenada iniciales, además del número de veces que desea correr el ciclo para encontrar la pendiente y la ordenada.

$$printf("Ingrese el valor de la pendiente y el valor de las ordenadas iniciales"); \quad (29)$$

Aquí se guarda el dato ya formateado del stdin en la variable m0.

$$scanf("%f", m0); \quad (30)$$

Se abre el archivo donde se escribirán los resultados obtenidos del cálculo:

$$resultados = fopen("resultados", "w"); \quad (31)$$

Se inicia un ciclo donde el contador *j* inicia en cero, la condición de seguir realizando el ciclo es que se repita un número menor al valor de la variable "delta" y que al contador se le aumente una unidad cada vez que se completa la serie de instrucciones. Dicho ciclo encontrará la pendiente y la ordenada el número de veces que le haya indicado el usuario.

$$for(j = 0; j < delta; j++) \quad (32)$$

Se generan números al azar entre -1,1 y se le restan a los valores iniciales guardados en m0,b0 respectivamente. Dicho valor se guarda en la variable mn,b0 respectivamente. El signo "*" indica una multiplicación y el signo "/" una división.

$$mn = m0 - (1) + (2) * (rand()) / (double)RAND_MAX; \quad (33)$$

Se llaman y se hace uso de las funciones percentx y percenty declaradas y descritas en el programa funcion.c y los valores de los cálculos de dichas funciones se guardan en px0, pxn y pxn2.

$$px0 = percentx(m0, b0); \quad (34)$$

$$pxn = percentx(m0, b0); \quad (35)$$

$$pxn2 = percenty(m0, b0); \quad (36)$$

Se crea una condición if donde si el valor guardado en pxn es mayor al valor guardado en px0 se igualaran los nuevos valores de mn y bn en m0 y b0 respectivamente.

$$if(px0 < pxn)\{ \quad (37)$$

$$m0 = mn; b0 = bn \quad (38)$$

else indica que si las condiciones anteriores dentro del if no se cumplen, se buscará un nuevo número aleatorio, y si dicho número es menor que pxn2 (alejpxn2) entonces dicho valor nuevo mn,bn se guardará en las variables m0,b0. RAND_MAX es una constante que representa el valor más grande al azar que una función puede regresar.

$$else\{ \quad (39)$$

$$ale = rand() / (double)RAND_MAX; \quad (40)$$

$$if(ale < pxn2) \quad (41)$$

$$m0 = mn; b0 = bn \quad (42)$$

Con fprintf se imprimen los valores de m0,b0 obtenidos a cada término de ciclo, y resultados indica que en dicho archivo se imprimirán.

$$fprintf(resultados, "%f %f ", m0, b0); \quad (43)$$

Se cierra el archivo de resultados

```
fclose(resultados); (44)
```

Indica si la secuencia de instrucciones sucedio correctamente, de lo contrario enviará signo de error

```
return 0; (45)
```

Makefile

El Makefile nos ayuda compilar un programa desde el origen del código.

Se incluyen los nombres de los programas que se desean compilar simultáneamente.

```
mainp : func.c mainp.c libs (46)
```

Gcc compila los programas func.c y mainp.c, junto con la librería libs en un programa llamado mainp.out, el -lm indica que se usó la librería con funciones matemáticas

```
gcc func.c mainp.c -o mainp.out -lm -I ./libs (47)
```

Clean limpia el ejecutable, que es el programa mainp.out

```
clean : (48)
```

```
rm mainp.out (49)
```

Cuando se coloca la palabra "make" en la terminal, dentro de la carpeta donde se encuentran los programas, librerías y documentos que forman parte del ejecutable; y se presiona enter "mainp.out" se compilará y ejecutará.

graph1.ipynb

Este programa se usó para graficar los datos obtenidos de los resultados del programa ejecutado con Makefile, llamado datos.txt, mientras se encuentre dividido en 2 columnas de datos numericos.

Se llaman las librerías que se desean usar, en este caso serán las estándar y para funciones matemáticas.

```
import matplotlib.pyplot as plt (50)
```

```
import numpy as np (51)
```

Se leen los datos del archivo "resultados", se grafica con matplotlib y se guarda la gráfica con plt.savefig.

```
file = "resultados" (52)
```

$$data_img = np.loadtxt(file) \quad (53)$$

$$plt.savefig('frecuencias absoluta pendientes.png') \quad (54)$$

Se leen los datos con los comandos mencionados anteriormente y se guardan en dos arreglos.

$$plt.savefig('frecuencias absoluta pendientes.png') \quad (55)$$

$$a = data_img[:, 0] \quad (56)$$

$$b = data_img[:, 1] \quad (57)$$

$$x = np.linspace(-5, 5, 100) \quad (58)$$

Se grafican en un rango de -5 a 5.

$$for i in range(1, 1000): \quad (59)$$

$$plt.plot(x, a[i * 10] * x + b[i * 10]) \quad (60)$$

Se grafican como pendiente y eje respectivamente y los resultados se guardan en archivos ".png"

$$for i in range(1, 1001) \quad (61)$$

$$file = name[i] \quad (62)$$

$$plt.plot(x, a[i * 10] * x + b[i * 10]) \quad (63)$$

$$plt.savefig(file + '.png') \quad (64)$$

3 Resultados

El programa se ejecutó con éxito y realizó el proceso de buscar la pendiente y la ordenada en "y" que más concordaba con el set de 10 de coordenadas (x,y) dadas. Dicho proceso se realizó un número de 10000 veces, que fué el valor de delta dado por el usuario. Las gráficas realizadas fueron las siguientes:

Los datos obtenidos en python nos muestran la tendencia de los datos obtenidos por el programa. En seguida se muestran los datos reales de la pendiente m y la ordenada b para comparar que tan cerca estuvo el programa de calcular los valores reales buscados:

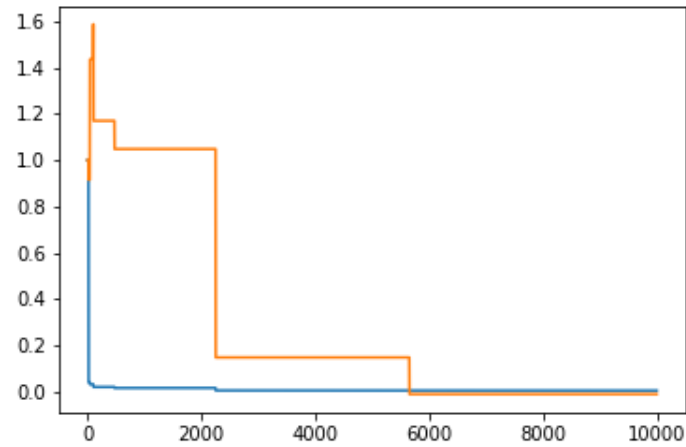


Figure 3: Se observa la tendencia de los valores de la pendiente y la ordenada.

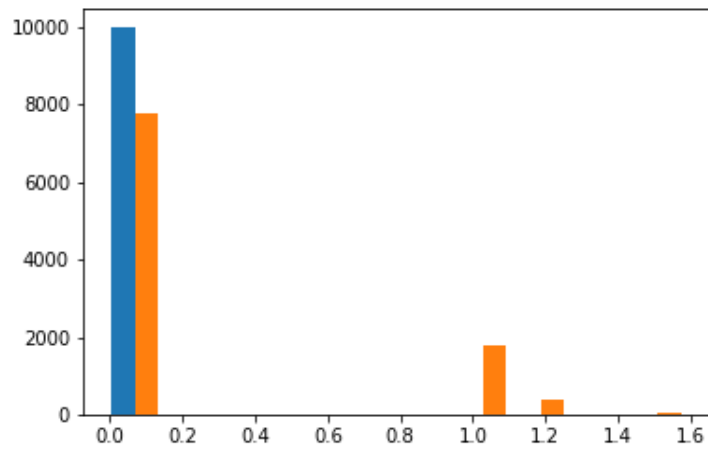


Figure 4: Histograma que representa el número de veces que se buscó el valor que mas se acercaba a la pendiente de los datos (x,y) ingresados y el valor de la misma.

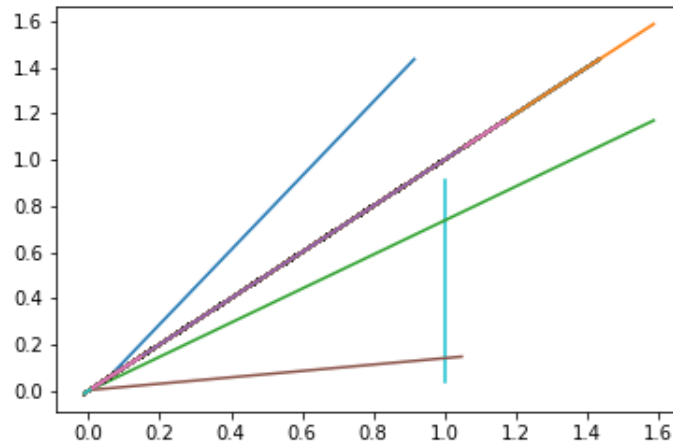


Figure 5: Pendientes encontradas.

10000 0.003527 -0.011711

Figure 6: Se observan: el número de veces que se buscó la pendiente y la ordenada, y el valor de la pendiente y valor de la ordenada encontrados.

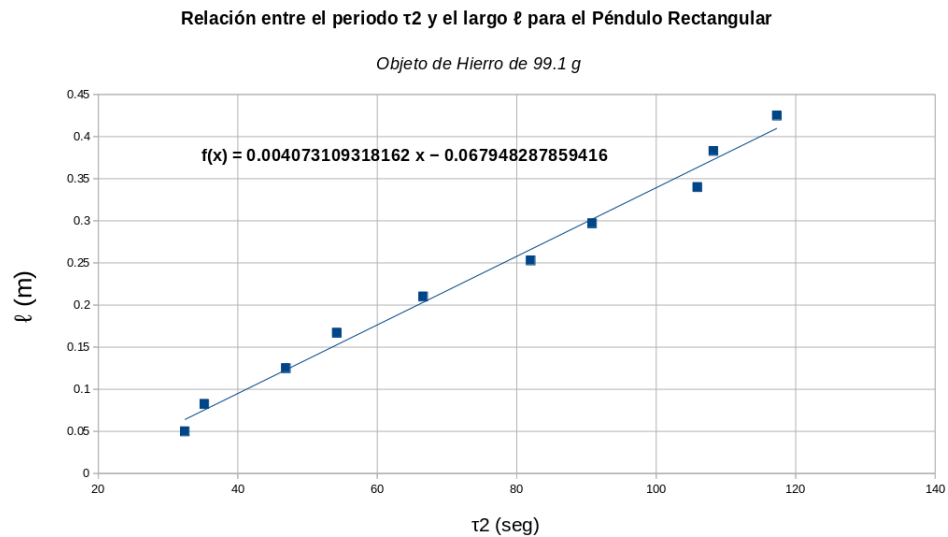


Figure 7: Se observan los puntos (x,y) con los cuales se ejecutó el programa y los valores reales de la pendiente m y la ordenada b .

4 Conclusiones

El programa usó el teorema de Bayes para encontrar la probabilidad de que el valor la pendiente de una gráfica se ajustara a los datos de 10 coordenadas dadas en x, y respectivamene. Al final, se notó la gran similitud entre los valores de los resultados obtenidos de la pendiente m y la ordenada b con el programa realizado, comparado con el valor real de la pendiente m y la ordenada b . Dichos resultados también se pueden obtener con el método de mínimos cuadrados, sin embargo, se escogió este programa debido a que el teorema de Bayes es una de las herramientas que ayuda a llegar al porcentaje de veracidad de un dato tomando en cuenta información ya conocida, también llamada información "a priori". Dicha herramienta tiene varias aplicaciones no solo para el área de la ciencia, también para el área de estadística y economía. El teorema de Bayes fué un tesoro que él dejó enterrado y que por desgracia, fué descubierto hasta después de su muerte. Ahora nos es de gran ayuda en sus aplicaciones para el uso y manejo de datos.