

402 MÉTODO DE DESCENSO DE GRADIENTE - TAREA 1

Octubre 06, 2025

Pregunta 1

La función de Himmelblau es un estándar para probar algoritmos de optimización, ya que posee cuatro mínimos locales idénticos. El objetivo es implementar y comparar el rendimiento del método de descenso de gradiente (GD) estándar con el descenso de gradiente con momento (GDM). Se debe analizar cómo la elección de los hiperparámetros (tasa de aprendizaje y momento) afecta la convergencia.

La función está definida como:

$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$$

Actividades:

1. Implementación: Defina la función de Himmelblau y su gradiente. Implemente el algoritmo de GD y GDM.
2. Experimentación:
 - Ejecute ambos algoritmos desde el punto inicial $(-2.5, 2.5)$.
 - Para GD, pruebe con tasas de aprendizaje $\alpha = [0.001, 0.005, 0.01]$.
 - Para GDM, utilice una tasa de aprendizaje fija $\eta = 0.005$ y pruebe con coeficientes de momentos $\beta = [0.5, 0.7, 0.9]$.
 - Utilice un máximo de 10000 iteraciones y una tolerancia de $1e-6$ como criterios de detención.

Código Fuente de la Implementación

```
from numpy import array as B, linspace as L, meshgrid as M
from matplotlib import pyplot as P

# Vectores
e = [B([-2.5, 2.5]), B([-2.5, 2.5]), B([-2.5, 2.5])]
ε = [B([-2.5, 2.5]), B([-2.5, 2.5]), B([-2.5, 2.5])]

# Tasas de aprendizaje
α = [0.001, 0.005, 0.01]
η = 0.005

# Factor de momento
β = [0.5, 0.7, 0.9]

# Velocidad
ω = [B([0, 0]), B([0, 0]), B([0, 0])]

#  $f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$ 
f = lambda x: (x[0]**2 + x[1] - 11)**2 + (x[0] + x[1]**2 - 7)**2
#  $\nabla f(x, y) = [4x(x^2 + y - 11) + 2(x + y^2 - 7), 2(x^2 + y - 11) + 4y(x + y^2 - 7)]$ 
ϕ = lambda x: B(
    [
        4 * x[0] * (x[0]**2 + x[1] - 11) + 2 * (x[0] + x[1]**2 - 7),
        2 * (x[0]**2 + x[1] - 11) + 4 * x[1] * (x[0] + x[1]**2 - 7),
    ]
)

#  $v(x, y, \beta, \omega) = \beta\omega + (1 - \beta)\nabla f(x, y)$ 
v = lambda w, x, b: b * w + (1 - b) * ϕ(x)
#  $g(x, y, \alpha) = (x, y) + \alpha\nabla f(x, y)$ 
g = lambda x, a: x - a * ϕ(x)
#  $\gamma(x, y, \omega) = (x, y) - \eta\omega$ 
γ = lambda x, w: x - η * w

# Historial de vectores
y = [[], [], []]
u = [[], [], []]

# Costo
c = [[], [], []]
x = [[], [], []]

# Iteraciones
i = [0, 0, 0]
ι = [0, 0, 0]

# Inicialización
y[0].append(e[0].copy())
y[1].append(e[1].copy())
y[2].append(e[2].copy())
u[0].append(ε[0].copy())
u[1].append(ε[1].copy())
u[2].append(ε[2].copy())
c[0].append(f(e[0]))
c[1].append(f(e[1]))
c[2].append(f(e[2]))
x[0].append(f(ε[0]))
x[1].append(f(ε[1]))
x[2].append(f(ε[2]))

for _ in range(10000):
    e[0] = g(e[0], α[0])
    y[0].append(e[0].copy())
    c[0].append(f(e[0]))
    i[0] += 1
    if f(e[0]) < 1e-6:
        break

for _ in range(10000):
    e[1] = g(e[1], α[1])
    y[1].append(e[1].copy())
    c[1].append(f(e[1]))
    i[1] += 1
```

```

        if f(e[1]) < 1e-6:
            break

for _ in range(10000):
    e[2] = g(e[2], a[2])
    y[2].append(e[2].copy())
    c[2].append(f(e[2]))
    i[2] += 1
    if f(e[2]) < 1e-6:
        break

for _ in range(10000):
    w[0] = v(w[0], e[0], beta[0])
    e[0] = y(e[0], w[0])
    u[0].append(e[0].copy())
    k[0].append(f(e[0]))
    i[0] += 1
    if f(e[0]) < 1e-6:
        break

for _ in range(10000):
    w[1] = v(w[1], e[1], beta[1])
    e[1] = y(e[1], w[1])
    u[1].append(e[1].copy())
    k[1].append(f(e[1]))
    i[1] += 1
    if f(e[1]) < 1e-6:
        break

for _ in range(10000):
    w[2] = v(w[2], e[2], beta[2])
    e[2] = y(e[2], w[2])
    u[2].append(e[2].copy())
    k[2].append(f(e[2]))
    i[2] += 1
    if f(e[2]) < 1e-6:
        break

# Gráficos
theta = pi(-5, 5, 100)
phi = pi(5, -5, 100)
M, N = m(theta, phi)
n.figure(figsize=(12, 10))
n.contour(M, N, f(B([M, N])), cmap="viridis", alpha=0.6, levels=20)
n.contourf(M, N, f(B([M, N])), cmap="viridis", alpha=0.1, levels=20)
mu = B(
    [[3.0, 2.0], [-2.805118, 3.131312], [-3.779310, -3.283186], [3.584428, -1.848126]]
)
n.scatter(
    mu[:, 0],
    mu[:, 1],
    color="red",
    marker="*",
    s=200,
    label="Mínimos globales",
    zorder=5,
)

u = ["blue", "green", "purple", "orange", "brown", "pink"]
a = [f"DG  $\alpha = \{\xi\}$ " for xi in alpha] + [f"DGM  $\beta = \{\eta\}$ " for eta in beta]

for eta in range(0, 6, 1):
    if eta < 3:
        Y = B(y[eta])
        n.plot(
            Y[:, 0],
            Y[:, 1],
            color=u[eta],
            linewidth=2,
            label=a[eta],
            alpha=0.1 if eta != 2 else 1,
        )
        n.scatter(Y[0, 0], Y[0, 1], color=u[eta], marker="o", s=100, zorder=4)
        n.scatter(Y[-1, 0], Y[-1, 1], color=u[eta], marker="s", s=100, zorder=4)
    else:

```

```

Y = v(u[j] - 3))
n.plot(
    Y[:, 0],
    Y[:, 1],
    color=u[j],
    linewidth=2,
    label=v[j],
    linestyle="--",
    alpha=0.1 if j != 3 else 1,
)
n.scatter(Y[0, 0], Y[0, 1], color=u[j], marker="o", s=100, zorder=4)
n.scatter(Y[-1, 0], Y[-1, 1], color=u[j], marker="s", s=100, zorder=4)

n.xlabel("x")
n.ylabel("y")
n.title("Trayectorias de Convergencia en la Función de Himmelblau")
n.legend()
n.grid(True, alpha=0.3)
n.axis("equal")
n.xlim(-5, 5)
n.ylim(-5, 5)

# Tabla comparativa
print("\n" + "=" * 90)
print("TABLA COMPARATIVA: GRADIENTE DESCENDENTE vs GRADIENTE DESCENDENTE CON MOMENTO")
print("=" * 90)
print(
    f"{'Algoritmo':<15} {'Hiperparámetros':<20} {'Mínimo (x,y)':<25} {'f(x,y)':<12}"
    f"{'Iteraciones':<12}"
)
print("-" * 90)

# Resultados DG
for j in range(3):
    print(
        f"{'DG':<15} {'α={α[j]}':<20} ({float(e[j][0]):<10.6f}, {float(e[j][1]):<11.6f})"
        f"{'f(e[j]):<12.3e} {'i[j]:<12}'"
    )

# Resultados DGM
for j in range(3):
    print(
        f"{'DGM':<15} {'α={α[j]}, β={β[j]}':<20} ({float(ε[j][0]):<10.6f}, {float(ε[j]"
        f"[1]):<11.6f}) {'f(ε[j]):<12.2e} {'v[j]:<12}'"
    )

print("=" * 90)

n.show()

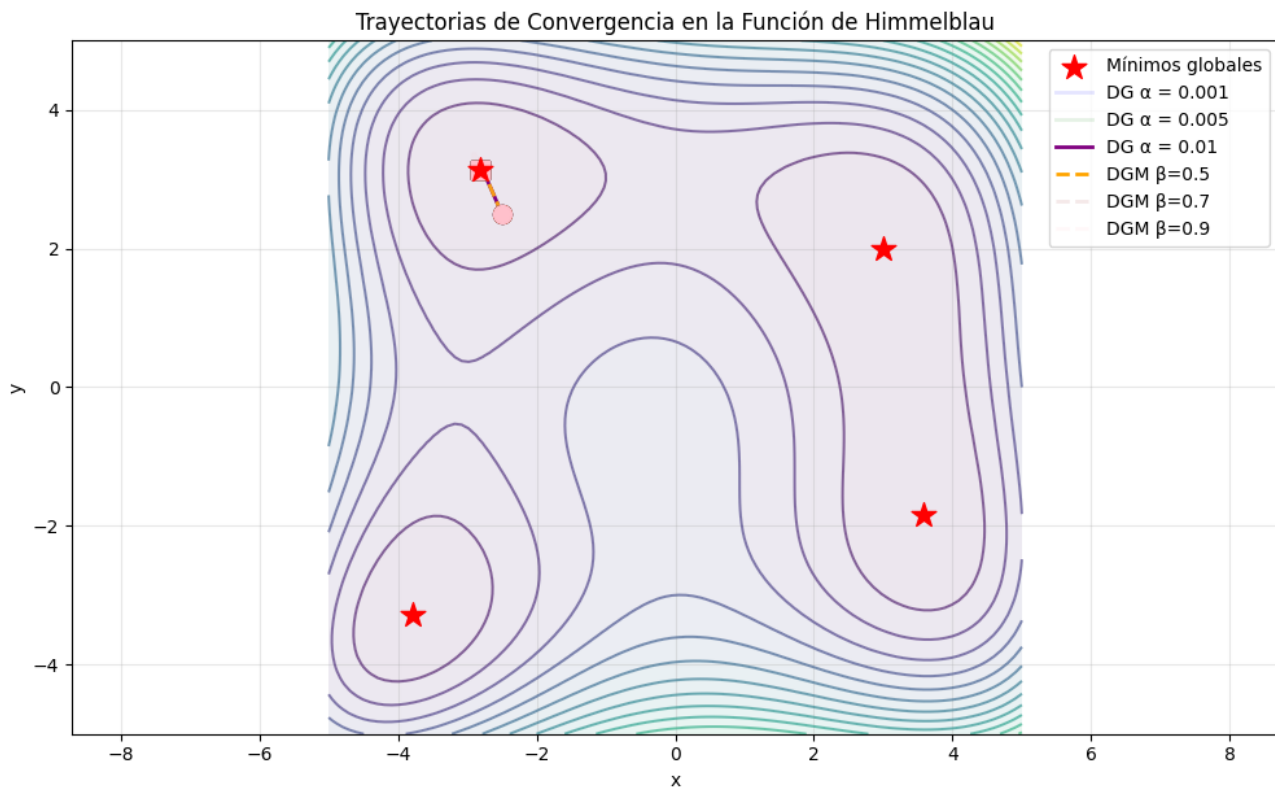
```

3. Análisis y visualización:

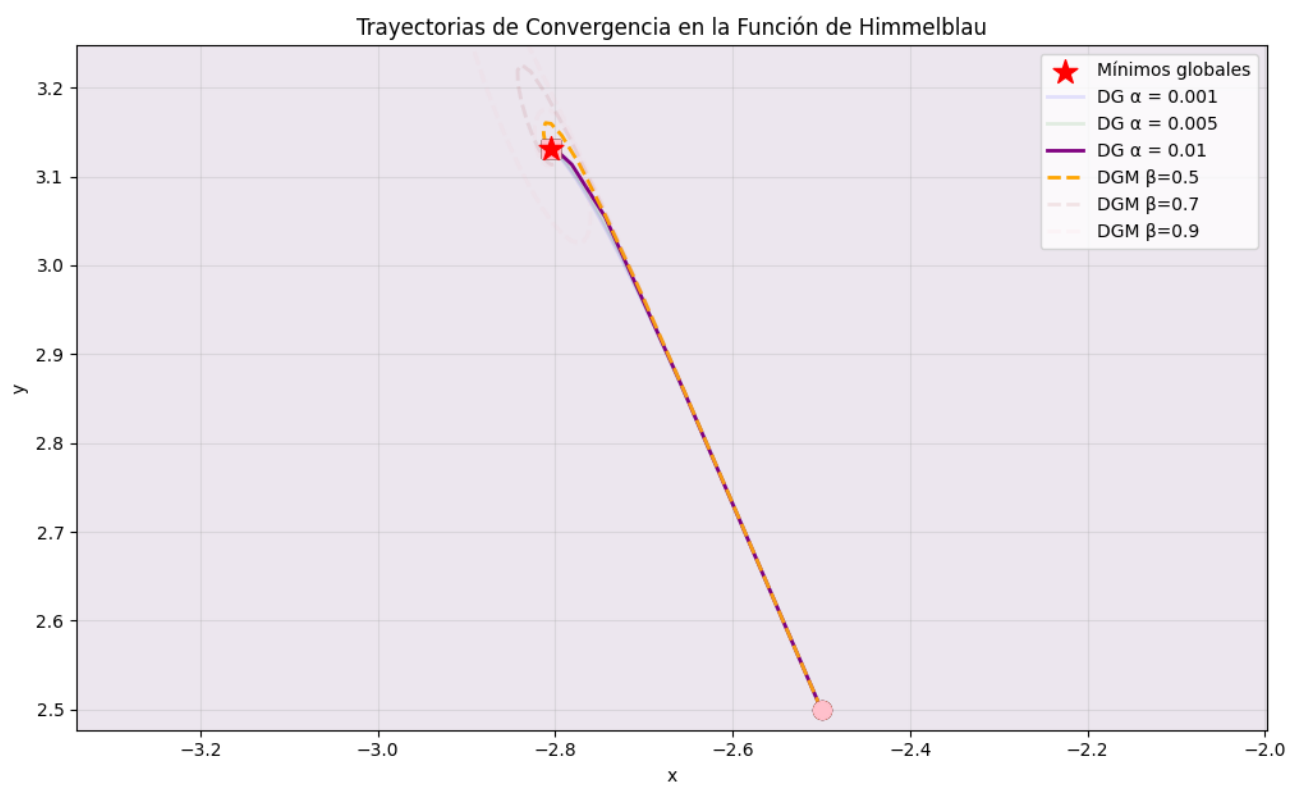
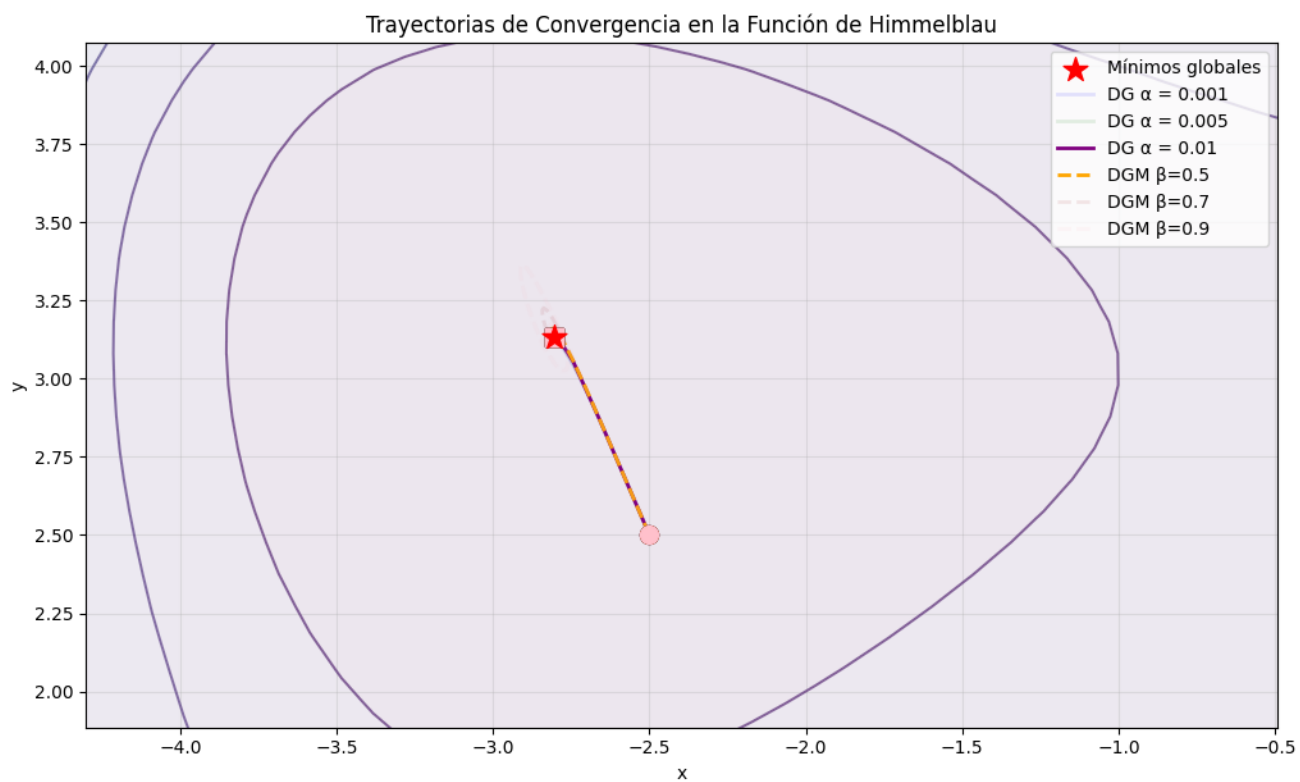
- Genere una tabla comparativa que muestre para cada configuración: el algoritmo, los hiperparámetros usados, el mínimo encontrado, el valor de la función en ese mínimo y el número de iteraciones requeridas.

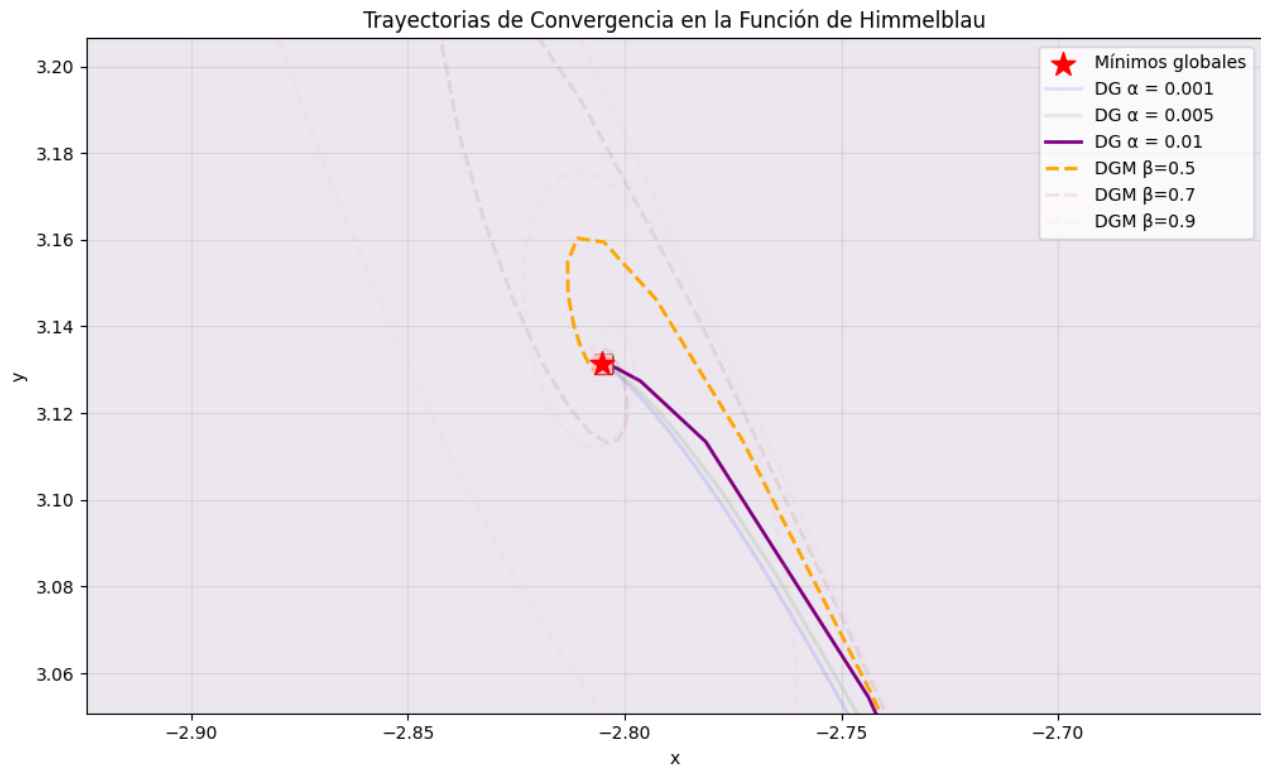
TABLA COMPARATIVA: GRADIENTE DESCENDENTE vs GRADIENTE DESCENDENTE CON MOMENTO				
Algoritmo	Hiperparámetros	Mínimo (x,y)	f(x,y)	Iteraciones
DG	$\alpha=0.001$	(-2.804961 , 3.131256)	9.214e-07	118
DG	$\alpha=0.005$	(-2.804993 , 3.131282)	5.411e-07	21
DG	$\alpha=0.01$	(-2.804978 , 3.131296)	6.434e-07	8
DGM	$\alpha=0.005, \beta=0.5$	(-2.805062 , 3.131431)	6.77e-07	25
DGM	$\alpha=0.005, \beta=0.7$	(-2.804985 , 3.131215)	9.45e-07	43
DGM	$\alpha=0.005, \beta=0.9$	(-2.805105 , 3.131250)	1.64e-07	135

- Genere un gráfico de contorno de la función de Himmelblau. Sobre este, grafique las trayectorias de convergencia para la mejor configuración de GD y la mejor configuración de GDM.



Aunque se pide sólo la graficación de los mejores resultados de cada conjunto de algoritmos, de igual manera se graficaron (con baja transparencia) los otros caminos.





- Escriba una conclusión breve explicando qué algoritmo convergió más eficientemente y por qué el momento ayuda (o perjudica) en este caso particular

El algoritmo que convergió más eficientemente fue el descenso del gradiente tradicional con la tasa de aprendizaje igual a 0.01, la cual consiguió caer por debajo de la tolerancia establecida en sólo 8 iteraciones, y el motivo es relativamente simple: el valor de α alto permite descender el gradiente de manera rápida sin pasarse. Directo al objetivo, la curva incluso se puede notar algo cruda, con giros puntiagudos debido a que realizó muy pocos saltos a diferencia de los otros descensos tradicionales.

Por otro lado, para este caso particular, el momento es un incordio, pues el momento causa que la curva en todos los casos falle la caída hacia el mínimo, lo cual sería práctico en funciones con mínimos locales para perseguir el mínimo global, aquí causa que el algoritmo de vueltas de manera innecesaria en torno al verdadero máximo. Un mejor resultado se obtiene al reducir el factor de momento, lo cual es lamentable, porque eso significa realizar igualmente los cálculos extra para momento cuando estos realmente aportan poco. Para la función de Himmelblau, el descenso tradicional con la tasa de aprendizaje de 0.01 es superior.

Pregunta 2

Una tasa de aprendizaje fija no siempre es óptima. Las estrategias de decaimiento pueden acelerar la convergencia y evitar oscilaciones cerca del mínimo. El objetivo, es implementar y comparar tres estrategias diferentes de decaimiento de la tasa de aprendizaje para minizar la función de Rosenbrock:

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$

Actividades:

1. Implementación:

- Modifique el algoritmo de descenso de gradiente para incorporar un planificador de tasa de aprendizaje.
- Implemente las siguientes tres estrategias de decaimiento para α_k en la iteración k :
 1. Decaimiento por tiempo: $\alpha_k = \alpha_0 / 1+k \cdot \tau$
 2. Decaimiento por pasos: $\alpha_k = \alpha_0 \cdot d^{\lfloor k/s \rfloor}$, donde d es un factor de decaimiento y s es el tamaño del paso
 3. Decaimiento exponencial: $\alpha_k = \alpha_0 \cdot e^{(-k \cdot \tau)}$

2. Experimentación:

- Minimice la función de Rosenbrock partiendo desde $(-1, -1)$.
- Utilice los siguientes parámetros: $\alpha_0 = 0.002$, $\max_iter = 10000$.
- Para el decaimiento por tiempo y exponencial, use: $\tau = \alpha_0 / 10000$.
- Para el decaimiento por pasos, use un factor $d = 0.5$ y un tamaño de paso $s = 2000$.
- Compare estos tres métodos con el descenso de gradiente con α fijo de 0.002

Código Fuente de la Implementación

```
from numpy import array as B
from matplotlib import pyplot as p
from math import floor as ϕ, exp as e

e = [B([-1, -1]), B([-1, -1]), B([-1, -1]), B([-1, -1])]
α = [0.002, 0.002, 0.002, 0.002]
d = 0.5
s = 2000
τ = α[0] / 10_000
i = [0, 0, 0, 0]
h = [[], [], [], []]
η = [[], [], [], []]

# f(x, y) = (1 - x)2 + 100(y - x2)2
f = lambda x: (1 - x[0]) ** 2 + 100 * (x[1] - x[0] ** 2) ** 2
# ∇f(x, y) = [ -2(1 - x) - 400x(y - x2), 200(y - x2) ]
ϕ = lambda x: B(
    [
        -2 * (1 - x[0]) - 400 * x[0] * (x[1] - x[0] ** 2),
        200 * (x[1] - x[0] ** 2),
    ]
)

# g(x, y, α) = (x, y) + α∇f(x, y)
g = lambda x, a: x - a * ϕ(x)

# Inicializar
h[0].append(f(e[0]))
h[1].append(f(e[1]))
h[2].append(f(e[2]))
h[3].append(f(e[3]))
η[0].append(α[0])
η[1].append(α[1])
η[2].append(α[2])
η[3].append(α[3])

for k in range(0, 10000, 1):
    e[0] = g(e[0], α[0])
    h[0].append(f(e[0]))
    η[0].append(α[0])
    i[0] += 1
    if f(e[0]) < 1e-6:
        break

for k in range(0, 10000, 1):
    α[1] = α[0] / (1 + k * τ)
    e[1] = g(e[1], α[1])
    h[1].append(f(e[1]))
    η[1].append(α[1])
    i[1] += 1
    if f(e[1]) < 1e-6:
        break

for k in range(0, 10000, 1):
    α[2] = α[0] * (d ** ϕ(k / s))
    e[2] = g(e[2], α[2])
    h[2].append(f(e[2]))
    η[2].append(α[2])
    i[2] += 1
    if f(e[2]) < 1e-6:
        break

for k in range(0, 10000, 1):
    α[3] = α[0] * e(-k * τ)
    e[3] = g(e[3], α[3])
    h[3].append(f(e[3]))
    η[3].append(α[3])
    i[3] += 1
    if f(e[3]) < 1e-6:
        break

F = p.figure(figsize=(15, 6))
```

```

A = n.subplot(1, 2, 1)
A.grid(True, linestyle="--", alpha=0.4)
A.set_title("α en función del número de iteraciones")
A.set_xlabel("Número de iteraciones (k)")
A.set_ylabel("Valor de α")

for j, (label, color) in enumerate(
    [
        ("α constante", "blue"),
        ("α por tiempo", "red"),
        ("α por pasos", "green"),
        ("α exponencial", "purple"),
    ]
):
    iterations = range(len(η[j]))
    A.plot(iterations, η[j], "-", label=label, color=color, linewidth=2)
A.legend()

B = n.subplot(1, 2, 2)
B.grid(True, linestyle="--", alpha=0.4)
B.set_title("Función de costo en función del número de iteraciones")
B.set_xlabel("Número de iteraciones (k)")
B.set_ylabel("Valor de la función de costo")
B.set_yscale("log") # Use log scale for better visualization

for j, (label, color) in enumerate(
    [
        ("α constante", "blue"),
        ("α por tiempo", "red"),
        ("α por pasos", "green"),
        ("α exponencial", "purple"),
    ]
):
    iterations = range(len(h[j]))
    B.plot(iterations, h[j], "-", label=label, color=color, linewidth=2)
B.legend()

n.tight_layout()

# Tabla comparativa
print("\n" + "=" * 90)
print("TABLA COMPARATIVA: DESCENSO DEL GRADIENTE ESTÁTICO vs DINÁMICO")
print("=" * 90)
print(
    f"{'Algoritmo':<15} {'Valores de α':<20} {'Mínimo (x,y)':<25} {'f(x,y)':<12}"
    f"{'Iteraciones':<12}"
)
print("-" * 90)

# Resultados
for j, k in enumerate(["α constante", "α por tiempo", "α por pasos", "α exponencial"]):
    print(
        f"{'k':<15} {'f'α = {α[j]:.6f}':<20} ({float(e[j][0]):<10.6f}, {float(e[j][1]):<11.6f})"
        f"{'f(e[j]):<12.3e} {'i[j]:<12}"
    )

print("=" * 90)

n.show()

```

3. Análisis y visualización:

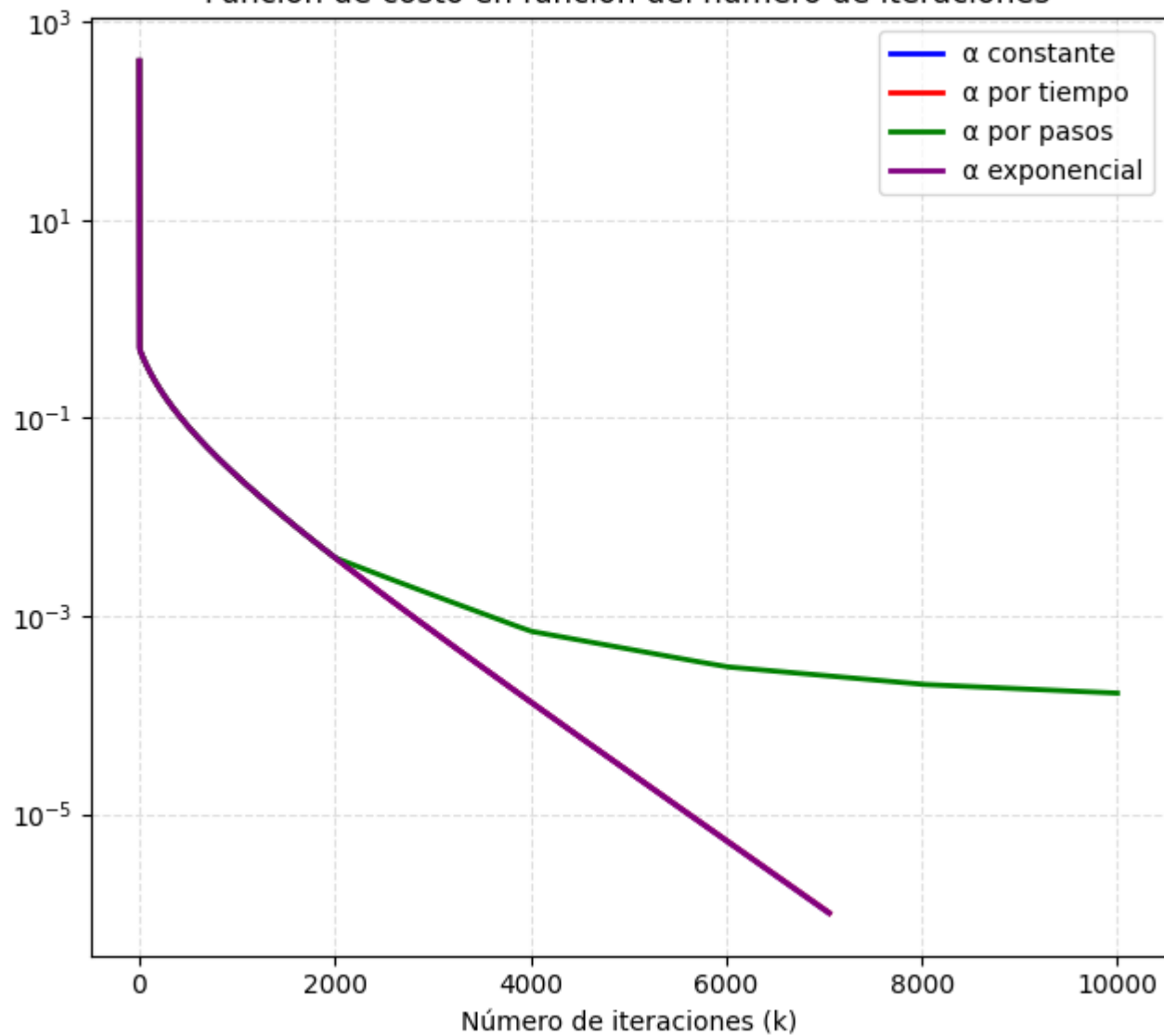
- Genere una tabla comparando las cuatro estrategias (incluida la de α fijo), mostrando el punto final alcanzado y el número de iteraciones necesarias para converger.

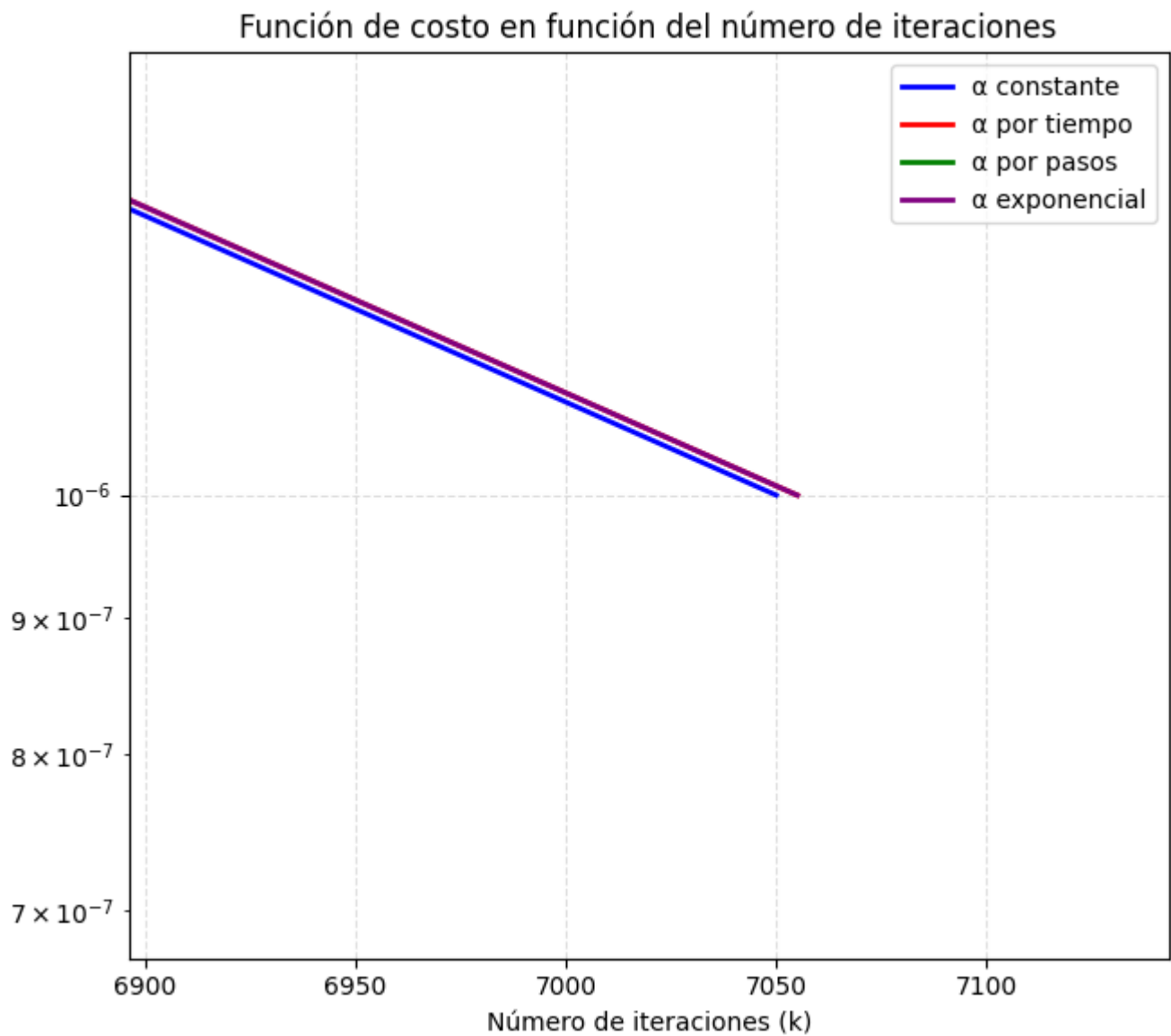
TABLA COMPARATIVA: DESCENSO DEL GRADIENTE ESTÁTICO vs DINÁMICO				
Algoritmo	Valores de α	Mínimo (x,y)	f(x,y)	Iteraciones
α constante	$\alpha = 0.002000$	(0.999001 , 0.997999)	9.998e-07	7050
α por tiempo	$\alpha = 0.001997$	(0.999001 , 0.997999)	9.997e-07	7055
α por pasos	$\alpha = 0.000125$	(0.987124 , 0.974362)	1.661e-04	10000
α exponencial	$\alpha = 0.001997$	(0.999001 , 0.997999)	9.997e-07	7055

Se puede apreciar cómo el algoritmo de actualización de la tasa de aprendizaje por pasos ni siquiera consigue converger, saliendo del bucle por acabarlo en lugar de optimizar el costo bajo la tolerancia elegida.

- Cree un gráfico que muestre la evolución del valor de la función objetivo $f(x, y)$ a lo largo de las iteraciones para cada una de las cuatro estrategias. Use una escala logarítmica en el eje Y para una mejor visualización.

Función de costo en función del número de iteraciones

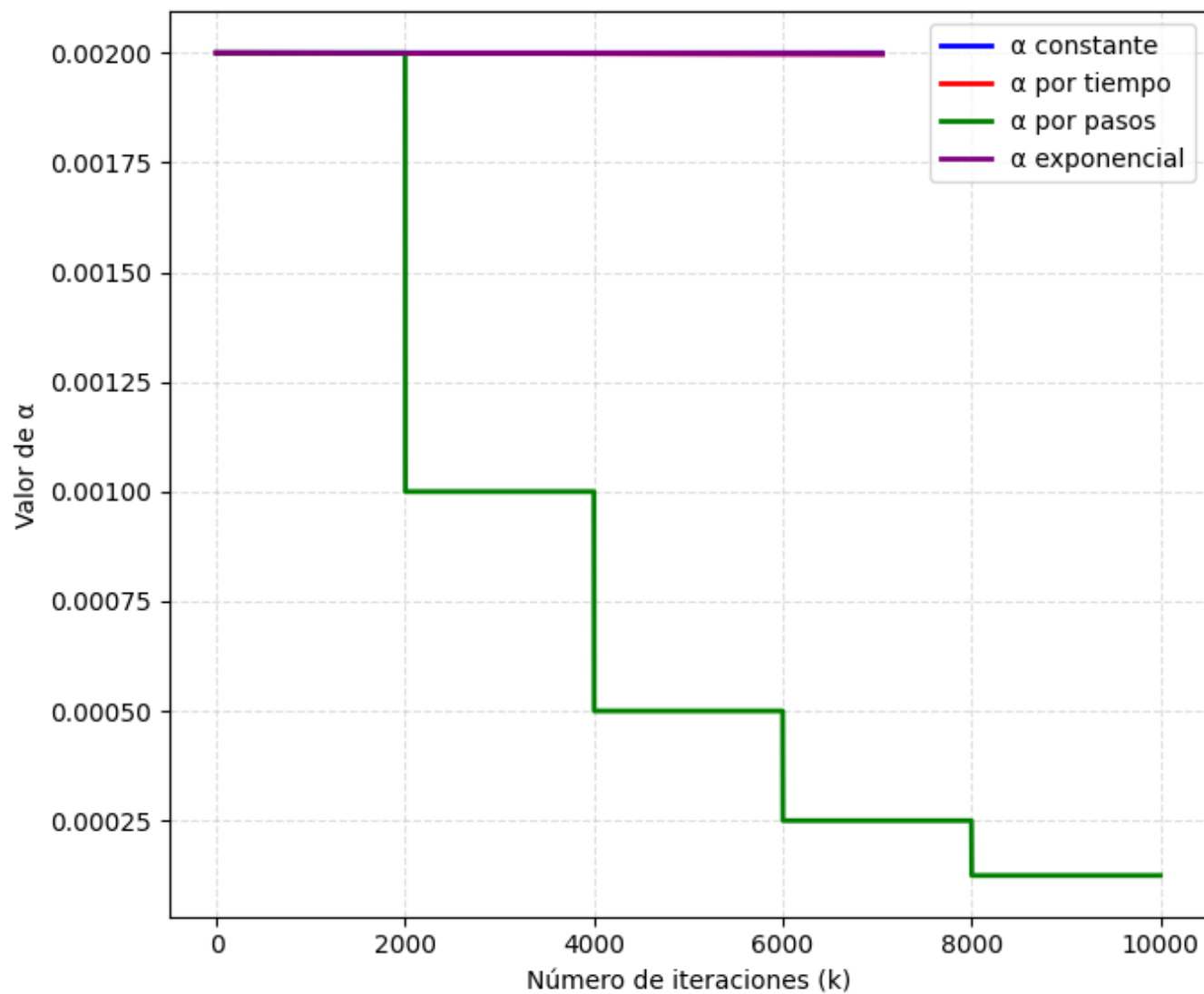


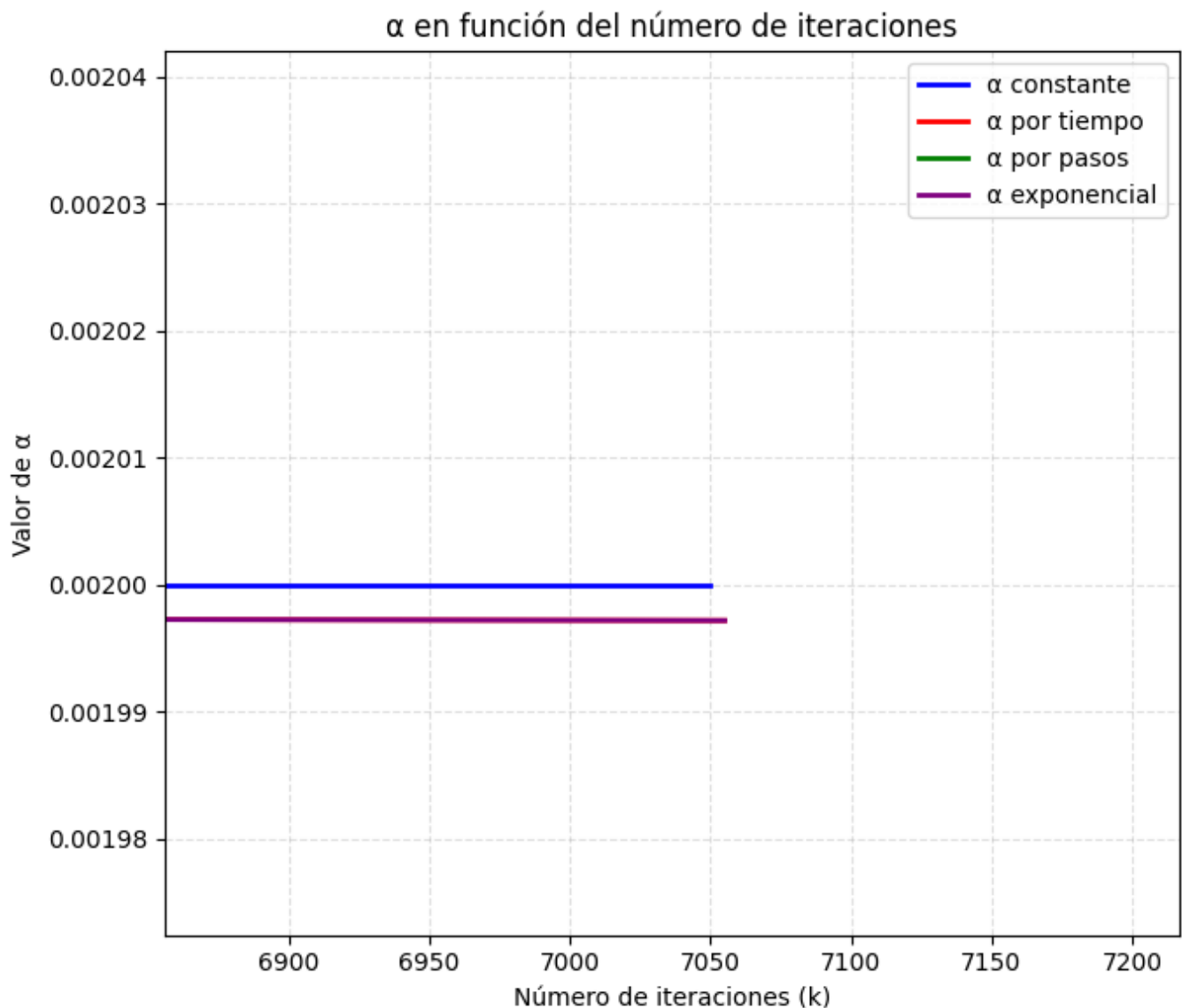


Zoom al gráfico del valor de la función de costo en función del número de iteraciones para poder visualizar cómo la tasa de aprendizaje constante consigue un mejor desempeño, además de que la tasa por tiempo y exponencial coinciden perfectamente, superponiéndose en el gráfico.

- En un segundo gráfico, muestre cómo cambia el valor de la tasa de aprendizaje α_k en cada iteración para las tres estrategias de decaimiento.

α en función del número de iteraciones





- Concluya cuál estrategia de decaimiento funcionó mejor para este problema y por qué.

La mejor estrategia resulta que es no utilizar ningún algoritmo de decaimiento, ¿Por qué? Pues porque tener un valor de α que decae significa mayor dificultad para moverse sobre la superficie tan plana de la función de Rosenbrock. Es por eso que sobre esta topología tan concreta, tener una tasa de aprendizaje constante proporciona un balance perfecto entre velocidad de convergencia y estabilidad, evitando la desaceleración prematura que arruina a algoritmos como el de pasos.

Aunque gráficamente pueda parecer que el algoritmo exponencial o decaimiento por tiempo tampoco tienen malos resultados, se debe considerar que estos algoritmos requieren cálculos extra para actualizar la tasa de aprendizaje, y por lo tanto se invierte más esfuerzo en calcular opciones inferiores a la tasa constante.

Pregunta 3

El método de mínimos cuadrados usado en la regresión lineal es muy sensible a los valores atípicos (outliers), ya que el error se eleva al cuadrado, magnificando el impacto de los puntos lejanos. Una alternativa más robusta es minimizar el Error Absoluto Medio (MAE).

La función objetivo MAE es:

$$f(\beta_0, \beta_1) = \frac{1}{m} \sum_{i=1}^m |y_i - (\beta_0 + \beta_1 x_i)|$$

Esta función no es diferenciable en todos sus puntos debido a la función de valor absoluto. Sin embargo, podemos usar su subgradiente para aplicar el descenso de gradiente. El subgradiente de la función objetivo es:

$$\frac{\partial f}{\partial \beta_j} = \frac{1}{m} \sum_{i=1}^m x_{ij} \cdot \text{sign}(y_i - \hat{y}_i)$$

donde $\text{sign}(z)$ es -1 si $z < 0$, $+1$ si $z > 0$, y 0 si $z = 0$.

El objetivo es implementar un algoritmo de regresión lineal robusto utilizando el descenso de subgradiente para minimizar el MAE y compararlo con el modelo estándar de Mínimos Cuadrados en presencia de datos atípicos.

Actividades:

1. Datos:
 - Utilice los datos de frecuencia_reloj y consumo_energia del apunte de clases.
 - Introduzca dos valores atípicos (outliers) en los datos para simular errores de medición, para la frecuencia de reloj [2.5, 4.5] y para el consumo de energia [130, 70].
2. Implementación:
 - Defina la función objetivo MAE y su subgradiente.
 - Implemente un algoritmo de descenso de gradiente que utilice el subgradiente MAE para encontrar los parámetros β óptimos.

Código Fuente de la Implementación

```
from numpy import (
    array as B,
    transpose as T,
    sum as d,
    sign as s,
    column_stack as K,
    ones as o,
    linspace as n,
)
from matplotlib import pyplot as p

f = B(
    [
        4.29,
        4.34,
        2.72,
        4.08,
        3.60,
        3.30,
        3.84,
        2.34,
        3.64,
        3.76,
        3.14,
        3.80,
        4.34,
        2.64,
        3.16,
        4.35,
        4.45,
        2.29,
        3.19,
        3.40,
        4.26,
        2.35,
        4.47,
        4.37,
        2.21,
        2.5,
        4.5,
    ]
)
f = K([o(f.shape[0]), f[:, None]])
c = B(
    [
        122,
        130,
        88,
        121,
        108,
        102,
        107,
        75,
        105,
        119,
        102,
        112,
        125,
        79,
        98,
        127,
        121,
        73,
        105,
        101,
        117,
        78,
        123,
        119,
        70,
        130,
        70,
    ]
)
```

```

    ]
)

ψ = B([2.5, 4.5])
ψ = κ([ö(ψ.shape[0]), ψ[:, None]])
ψ = B([130, 70])

θ = [B([0, 0]), B([0, 0])]
α = 0.01

h = lambda θ: f @ θ
J = lambda θ: д((h(θ) - c) ** 2) / (2 * f.shape[0])
j = lambda θ: (τ(f) @ (h(θ) - c)) / f.shape[0]
Й = lambda θ: д(abs(c - h(θ))) / f.shape[0]
й = lambda θ: (τ(f) @ ш(h(θ) - c)) / f.shape[0]

l = π(f[:, 1].min(), f[:, 1].max(), 100)

for _ in range(100000):
    θ[0] = θ[0] - α * j(θ[0])

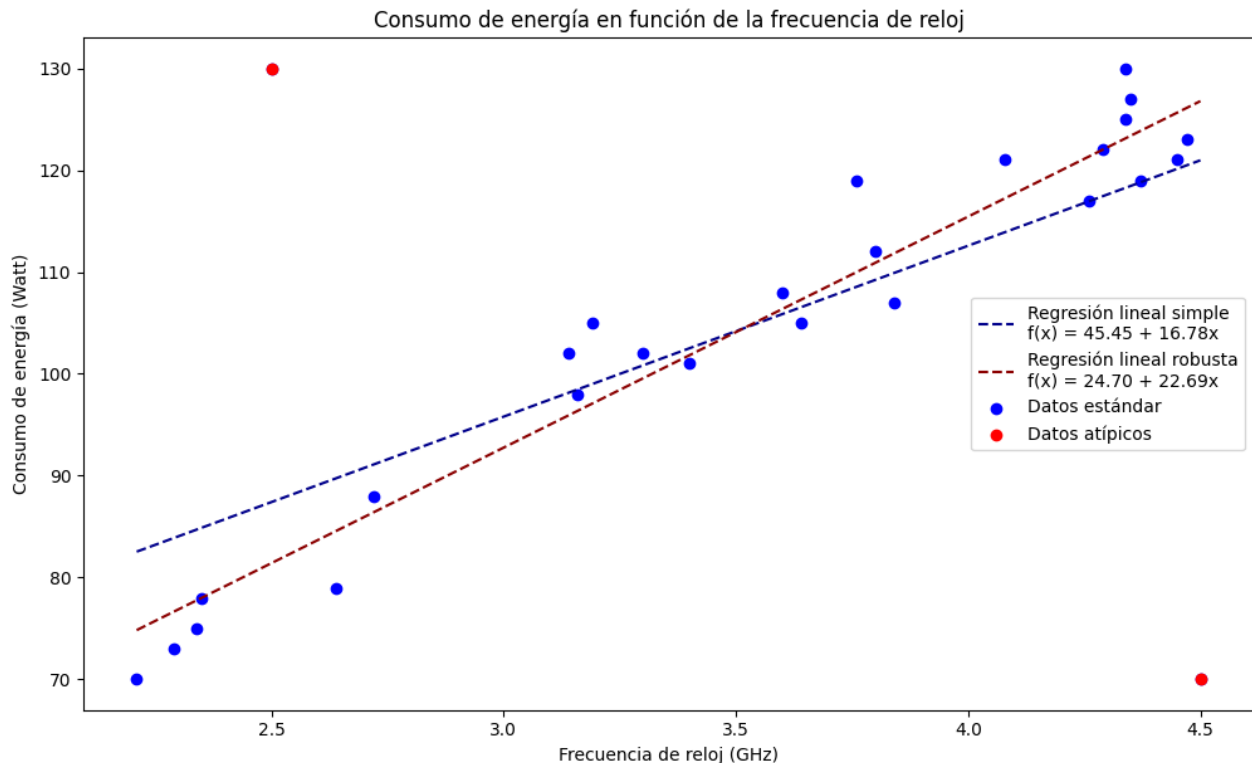
for _ in range(100000):
    θ[1] = θ[1] - α * й(θ[1])

π.figure(figsize=(12, 8))
π.plot(
    1,
    θ[0][1] * 1 + θ[0][0],
    color="darkblue",
    linestyle="--",
    label=f"Regresión lineal simple\nf(x) = {θ[0][0]:.2f} + {θ[0][1]:.2f}x",
)
π.plot(
    1,
    θ[1][1] * 1 + θ[1][0],
    color="darkred",
    linestyle="--",
    label=f"Regresión lineal robusta\nf(x) = {θ[1][0]:.2f} + {θ[1][1]:.2f}x",
)
π.scatter(f[:, 1], c, color="blue", marker="o", label="Datos estándar")
π.scatter(ψ[:, 1], ψ, color="red", marker="o", label="Datos atípicos")
π.xlabel("Frecuencia de reloj (GHz)")
π.ylabel("Consumo de energía (Watt)")
π.title("Consumo de energía en función de la frecuencia de reloj")
π.tick_params()
π.legend()
π.show()

```

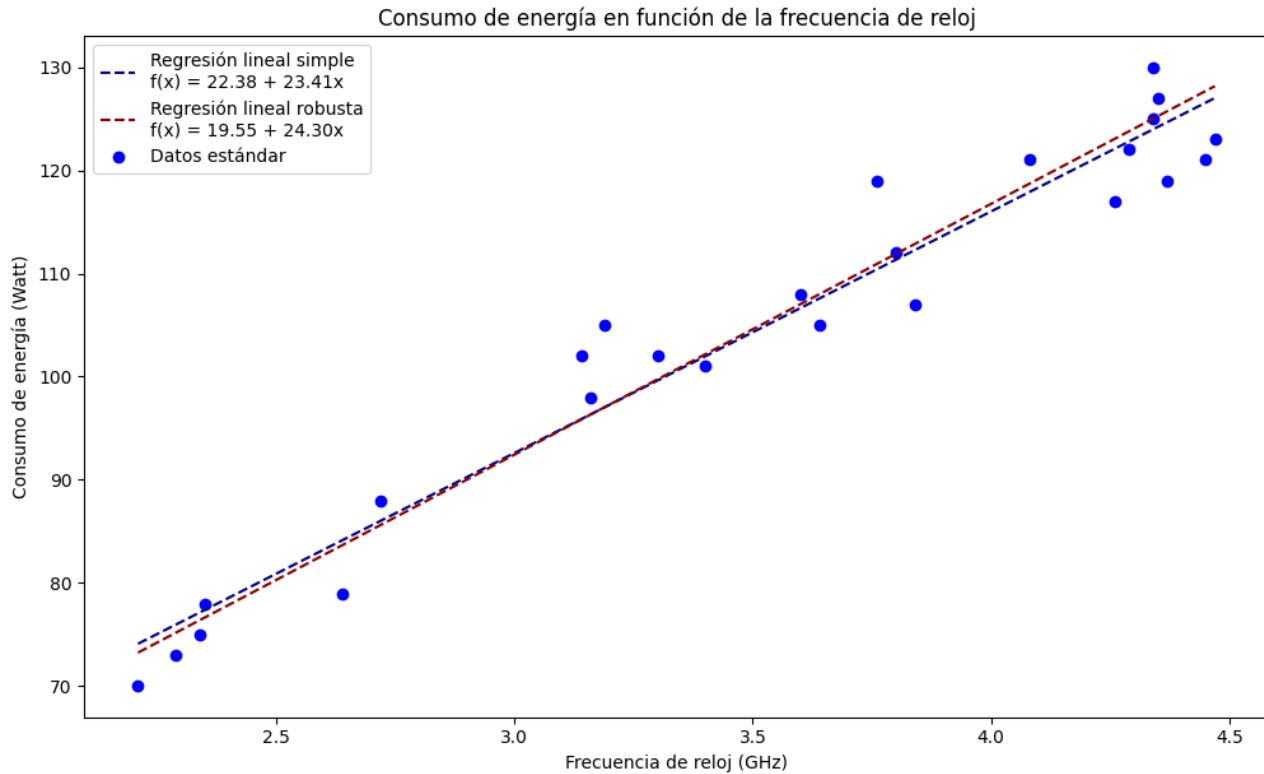
3. Experimentación y análisis:

- Ejecute el algoritmo de regresión lineal original (basado en MSE) sobre el conjunto de datos con outliers.
- Ejecute el nuevo algoritmo de regresión robusta (basado en MAE) sobre el mismo conjunto de datos con outliers.
- Gráfico: Cree un único diagrama de dispersión que muestre los datos con los outliers. Sobre este gráfico, dibuje las dos rectas de regresión obtenidas:
 1. La recta del modelo de Mínimos Cuadrados (MSE).
 2. La recta del modelo de Error Absoluto Medio (MAE).



- Análisis: En una breve conclusión, explique cuál de las dos rectas se ajusta mejor a la tendencia general de los datos originales. ¿Por qué el modelo basado en MAE es considerado más “robusto” ante valores atípicos? ¿Cómo se relaciona esto con la penalización que cada función de objetivo impone a los errores grandes?

El ejercicio nos pide trabajar en todo momento *con* los datos atípicos, pero si para responder la pregunta “¿Cuál de las dos rectas se ajusta mejor a la tendencia general de los *datos originales*?”, entonces haciendo unas ligeras modificaciones al algoritmo para quitar los datos atípicos podemos encontrar el siguiente gráfico:



No es difícil apreciar cómo la regresión lineal robusta se adapta mejor al caso con datos atípicos, pero es que también obtiene mejor precisión que la regresión lineal simple en el caso *sin datos atípicos*, demostrando que en términos de robustés, sí, el algoritmo es notoriamente mejor en todo aspecto.

Los datos atípicos dentro del algoritmo simple se someten a una exponenciación, lo cual causa que los datos de mayor magnitud tengan un mayor peso a la hora desviar la función final. El algoritmo MAE tiene una gran ventaja aquí, jerarquizando los pesos al no utilizar el mecanismo de exponenciación, y en su lugar utilizando un valor absoluto (que realmente podría interpretarse como “calcular la raíz cuadrada de los datos exponenciados”, lo cual reduce nuevamente el boom que los datos atípicos obtienen en sus pesos) lo cual mantiene los datos más a raya, y convergiendo notoriamente más rápido a regresiones más precisas sin la necesidad de muchos más cálculos que una MSE de toda la vida.