	Homework #1 Answers Chongye Feng
In [1]: In [2]:	from mpl_toolkits.mplot3d import Axes3D import numpy as np %matplotlib inline import time
	<pre>def timeit(f): def timed(*args, **kw): ts = time.time() result = f(*args, **kw) te = time.time() print('func:%r took: %2.4f sec' % (fname, te-ts)) return result return timed</pre>
	 a: because it increases the chance of finding the minimum value of the function by evenly dividing the larger interval in half. This allows for a more thorough search of the potential minimum in comparison to only searching a portion of the interval. Additionally, bisecting the larger interval also reduces the overall size of the interval, making it easier to converge towards the minimum value. b: If point e is placed in [a,b]:
	1. if f(e) > f(b), we could reduce the interval length from 0.75 to 0.5, a 1/3 reduction; 2. if f(e) < f(b), we could reduce the interval length from 0.75 to 0.5, a 1/3 reduction; 3. The average reduction would be 1/3. • If point e is placed in [b,d]: 1. if f(e) > f(b), we could reduce the interval length from 0.75 to 0.625, a 1/6 reduction; 2. if f(e) < f(b), we could reduce the interval length from 0.75 to 0.25, a 1/2 reduction; 3. The average reduction would be 1/3. • c: after step 2, the new interval would be, if e is placed in [a,b], [a,e,b] or [e,b,d] in the size of 0.5. And both of the oiginal bisection point is in the center of the interval, which is in a similar situation of step 1. Take [a,e,b] as an example, both [a,e] and [e,b] are in the same size and no difference without searching: • if f is in [a,e], f=0.125; • if f is in [e,b], f=0.375; • after getting f, the new reduced interval would be in size 0.25 or 0.375; • d: 1/2
<pre>In [3]: Out[3]:</pre>	• e: # step 1 (0.5 + 0.75) / 2 0.625
<pre>In [4]: Out[4]: In [5]:</pre>	<pre># step 2 (0.5 + 0.75 + 4/3) / 4 0.64583333333333333 # step 3 ((0.5 + 0.75) * 3 + 4/3) / 8</pre>
Out[5]:	0.6354166666666666 • f: For Golden Section, the ratio of size of interval of the previous step is 0.618, no change. 0.618
In [6]:	<pre>def func_q2(coord): x = coord[0] y = coord[1] return x**4 - x**2 + y**2 + 2*x*y -2 x = np.linspace(-2, 2, 200) y = np.linspace(-2, 2, 200) x, Y = np.meshgrid(x, y) z = func_q2((X, Y)) fig = plt.figure() ax = fig.add_subplot(ill, projection='3d') ax.plot_surface(X, Y, Z, cmap=cm.coolwarm, linewidth=0, antialiased=False) ax.set_xlabel('X axis') ax.set_ylabel('Y axis') ax.set_plabel('Y axis') ax.set_plabel('Y axis') ax.view_init(45,45) plt.show()</pre>
	20 15 15 0 -2 -1 -2 -1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
<pre>In [7]: Out[7]: In [8]:</pre>	21 26601101 20 75022147 5 07102051
Out[8]: In [9]:	<pre>7.5625 def first_derivate_q2(coord): x = coord[0] y = coord[1] dx = 4 * x**3 - 2 * x + 2 * y dy = 2 * y + 2 * x return (dx, dy)</pre>
	first_derivate_q2((1.5,1.5)) (13.5, 6.0) So, the first step would be 0.1 * (13.5, 6), which is (1.35, 0.6).
<pre>In [11]: Out[11]:</pre>	And the x_0 (1.5, 1.5) with the limit is ([-2,2], [-2,2]). So, the x_1 would be (1.5 - 1.35, 1.5 - 0.6), which is (0.15, 0.90). To judge if it is a good step: $ func_q 2((0.15,0.9)) < func_q 2((1.5,1.5)) $ True $ f(x_1) < f(x_0) $ So that it is a good step.
	Based on this outcome, the step size would be increased to (0.12 * first_derivative(x,y)). And, the start point would be updated to x_1 (0.15, 0.90). • b: from pylab import * import numpy.linalg as LA
	<pre> timeit def steepest_descent(func, first_derivate, starting_point, stepsize, tol): # evaluate the gradient at starting point deriv = first_derivate(starting_point) count = 0 visited = [] while LA.norm(deriv) > tol and count < le6: # calculate new point position new_x = starting_point[0] - stepsize * deriv[0] new_y = starting_point[1] - stepsize * deriv[1] new_point = (new_x.new_y) visited.append(new_point) if func(new_point) < func(starting_point): # the step makes function evaluation lower - it is a good step. what do you do? starting_point = new_point deriv = first_derivate(starting_point) stepsize *= 1.2 else: # the step makes function evaluation higher - it is a bad step. what do you do? stepsize *= 0.5 </pre>
In [13]:	<pre>count += 1 # return the results return {"x": starting_point, "evaluation": func(starting_point), "n_steps": count} steepest_descent(func_q2,first_derivate_q2,(0.15,0.9),0.12,1e-5) func: 'steepest_descent' took: 0.0008 sec ('m', / 0.000000515461034, 0.0000060663318330)</pre>
Out[13]: In [14]:	<pre></pre>
In [15]:	<pre>@timeit def timming_BFGS(): res = minimize(func_q2, (1.5,1.5), method='BFGS', options={'disp': True, 'gtol': le-5}) print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_BFGS()</pre>
	Optimization terminated successfully. Current function value: -3.000000 Iterations: 7 Function evaluations: 24 Gradient evaluations: 8 Minimum of the function: -2.9999999999955 Location of the minimum: [0.99999979 -0.9999998] func: 'timming_BFGS' took: 0.0025 sec
In [16]:	<pre>@timeit def timming_CG(): res = minimize(func_q2, (1.5,1.5), method='CG', options={'disp': True, 'gtol': 1e-5}) print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_CG()</pre>
<pre>In [17]: Out[17]:</pre>	func: 'steepest_descent' took: 0.0006 sec
	'evaluation': -2.999999999999999999999999999999999999
In [19]:	<pre>x = coord[0] y = coord[1] return ((1 - x)**2 + 10 * (y - x**2)**2) def first_derivate_q3(coord): x = coord[0] y = coord[1] dx = 40 * x**3 + (2 - 40 * y) * x - 2 dy = -20 * x**2 + 20 * y</pre>
In [20]: Out[20]:	<pre>return (dx, dy) steepest_descent(func_q3,first_derivate_q3,(-0.5,1.5),0.1,1e-5) func:'steepest_descent' took: 0.0133 sec {'x': (0.9999908923749649, 0.9999815271830772), 'evaluation': 8.361266798228901e-11, 'n_steps': 1523}</pre>
In [21]:	<pre>b: **timeit def stochastic_gradient_descent(func.first_derivete, starting_point, stepsize, tol=1e-5, stochastic_injection=1): ""stochastic_injection: controls the magnitude of stochasticity (multiplied with stochastic_deriv) 0 for no stochasticity, equivalent to SD. Use 1 in this homework to run SG0 "evaluate the gradient at starting point deriv = first_derivate(starting_point) count=0 visited="["while IA.norm(deriv) > tol and count < le5: if stochastic_injection=0: # formulate a stochastic deriv that is the same norm as your gradient stochastic_deriv = normade.radn(len(starting_point)) stochastic_deriv = normade.radn(len(starting_point)) stochastic_deriv = stochastic_deriv / LA.norm(stochastic_deriv) * LA.norm(deriv) else:</pre>
In [22]:	<pre>else: # the step makes function evaluation higher - it is a bad step. what do you do? stepsize *= 0.5 count+=1 return {"x":starting_point,"evaluation":func(starting_point), "n_steps": len(visited)} stochastic_gradient_descent(func_g3,first_derivate_g3,(-0.5,1.5),0.1)</pre>
Out[22]:	func: 'stochastic_gradient_descent' took: 0.0551 sec
In [23]:	<pre>@timeit def timming_BFGS_q3(): res = minimize(func_q3, (-0.5,1.5), method='BFGS', options={'disp': True, 'gtol': 1e-5}) print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_BFGS_q3() Optimization terminated successfully. Current function value: 0.000000 Iterations: 22</pre>
In [24]:	Function evaluations: 93 Gradient evaluations: 31 Minimum of the function: 1.6856836004019217e-13 Location of the minimum: [0.99999959 0.99999917] func:'timming_BFGS_q3' took: 0.0064 sec @timeit def timming_CG_q3(): res = minimize(func_q3, (-0.5,1.5), method='CG', options={'disp': True, 'gtol': 1e-5})
	<pre>print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_CG_q3() Optimization terminated successfully.</pre>
	Gradient evaluations: 44 Minimum of the function: 2.0711814827200667e-13 Location of the minimum: [0.9999955 0.9999908] func: 'timming_G_g_q3' took: 0.0082 sec stochastic_gradient_descent(func_q3, first_derivate_q3, (-0.5, 1.5), 0.1) func: 'stochastic_gradient_descent' took: 0.0599 sec { 'x': (0.9999889255304705, 0.9999774304825385), 'evaluation': 1.2441376906426114e-10, 'n_steps': 2196} By comparing the three methods, BFGS and CG would take less steps than SGD method, 31 and 44 steps compared with 2369 steps regardly. So, in terms of step efficiency, BFGS method is better than CG method, steepest descent method is the worst in this case. • d: No, it is not possible to draw a firm conclusion on the outcome of optimization algorithms such as Stochastic Gradient Descent (SGD), Conjugate Gradient (CG), or BFGS with just one run of each method. To draw a robust conclusion on the outcome of different optimization algorithms, it is necessary to run multiple trials, with different initial conditions, such as random seed, tolerance, etc.
	I think it is better to perform a systematic and thorough evaluation of different optimization algorithms on a diverse range of test functions and real-world datasets, in order to obtain a comprehensive understanding of their performance and limitations. • e: steepest_descent(func_q3,first_derivate_q3,(-0.5,-1),0.1,1e-5)
Out[26]: In [27]: Out[27]:	func: 'steepest_descent' took: 0.0116 sec {'x': (0.99998565363), 0.9999786855407166), 'evaluation': 1.1081718190344482e-10, 'n_steps': 1479} stochastic_gradient_descent(func_q3,first_derivate_q3,(-0.5,1),0.1) func: 'stochastic_gradient_descent' took: 0.0529 sec {'x': (0.99998958181722, 0.9999788167573075), 'evaluation': 1.097153443375041e-10, 'n_steps': 2143} Non-Stochastic (Steepest Descent) method: n_steps
	1 1541 1503 Stochastic Gradient Descent method: n_steps -1 -0.5 -1 3224 2157 -0.5 2235 2305 0 3831 2263 0.5 2226 2176 1 2298 2027
	It shows that the non-stochastic (steepest descent) method would need less steps to converge, in general, than the Stochastic Gradient Descent method. At starting point (-1,0) and (-0.5,-0.5), because it is no "hill climbing" between starting point and local minimum point, the path is shorter comparing with the other points with the same x_coordinate, the steps needed is the lest among the points with same x_coordinate. It is observed in the non-stochastic (steepest descent) method. On the other hand, the Stochastic Gradient Descent method does not have this observation. Q_4 def func_g4(coord): x = coord[0] y = coord[1]
In [29]: In [30]:	<pre>return (2 * x**2 - 1.05 * x**4 + (1/6) * x**6 + x * y + y**2) def first_derivate_q4(coord): x = coord[0] y = coord[1] dx = x**5 - 4.2 * x**3 + 4 * x + y dy = x + 2 * y return (dx, dy)</pre> <pre> • a:</pre>
<pre>In [30]: Out[30]: In [31]:</pre>	<pre>stochastic_gradient_descent(func_q4,first_derivate_q4,(-1.5,-1.5),0.1) func:'stochastic_gradient_descent' took: 0.0020 sec {'x': (-1.7475528342658777, 0.8737809027052557), 'evaluation': 0.298638442258369, 'n_steps': 50} @timeit def timming_BFGS_q4(): res = minimize(func_q4, (-1.5,-1.5), method='BFGS', options={'disp': True, 'gtol': 1e-5}) print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_BFGS_q4()</pre>
In [32]:	Optimization terminated successfully. Current function value: 0.298638 Iterations: 8 Function evaluations: 30 Gradient evaluations: 10 Minimum of the function: 0.29863844223686 Location of the minimum: [-1.74755234 0.87377616] func: 'timming_BFGS_q4' took: 0.0027 sec @timeit def timming_CG_q4(): res = minimize(func_q4, (-1.5,-1.5), method='CG', options={'disp': True, 'gtol': le-5})
	<pre>print("Minimum of the function:", res.fun) print("Location of the minimum:", res.x) timming_CG_q4() Optimization terminated successfully.</pre>
	Comparing with 3e, the Stochastic Gradient Descent method required less steps in this case, 69 steps vs 2000+ steps. But the BFGS and CG methods have better performance than Stochastic Gradient Descent method, 8, 7 steps vs 69 steps. Some important observation, the Stochastic Gradient Descent method did not converge into the global minimum every time, and BFGS and CG methods converge into the same global minimum every time. It is possible for the Stochastic Gradient Descent method be traped by the local minimum. • b: @timeit def SGDM(func, first_derivate, starting_point, stepsize, momentum=0.9, tol=1e-5, stochastic_injection=1):
	<pre>visited = [] deriv = first_derivate(starting_point) previous_direction = np.zeros(len(starting_point)) while LA.norm(deriv) > tol and count < le5: deriv = first_derivate(starting_point) if stochastic_injection > 0: # formulate a stochastic_deriv that is the same norm as your gradient stochastic_deriv = np.random.normal(0, LA.norm(deriv), len(starting_point)) else: stochastic_deriv = np.zeros(len(starting_point)) direction = (deriv + stochastic_injection * stochastic_deriv) # use previous direction to compute the current direction direction = momentum * previous_direction + direction # calculate new point position new_point = starting_point - stepsize * direction if func(new_point) < func(starting_point): # the step makes function evaluation lower - it is a good step. starting_point = new_point</pre>
	<pre>starting_point = new_point previous_direction = direction stepsize *= 1.2 else: # the step makes function evaluation higher - it is a bad step. # if stepsize is too small, clear previous direction because we already know that is not a useful direction if stepsize < le-5: previous_direction = np.zeros(len(starting_point)) else: # decrease stepsize by factor of 2 and clear previous direction stepsize *= 0.5 previous_direction = np.zeros(len(starting_point)) visited.append(starting_point) count += 1 return {"x": starting_point, "evaluation": func(starting_point), "n_steps": len(visited)}</pre>
	SGDM(func_q4,first_derivate_q4,(-1.5,-1.5),0.1) func: 'SGDM' took: 0.0076 sec {'x': array([-1.74755189, 0.87377258]), 'evaluation': 0.29863844224940717, 'n_steps': 53} I don't get a better result using SGDM compared to SGD, CG or BFGS in finding the global minimum in terms of fewer steps. It took 0.0107 second and 78 steps to converge into the global minimum. I did several runs of the SGDM method. Sometimes it would be trapped by the local minimum, similar with SGD method. And the steps it takes on average is more than the SGD method, around 70 steps vs around 60 steps. In this Three-Hump Camel function, the rank of steps of methods to converge (from most to least) is SGDM > SGD > BFGS > CG.