```
Consider a multivariate function f: \mathbb{R}^n \to \mathbb{R} that
represents some cost or value. This is called an
objective function, and we often want to find an objective function
X \in \mathbb{R}^n that yields f's extremum: minimum or extremum
maximum, depending on whichever is
desirable.
It is important to note however that some
functions have no finite extremum. Other
                                                        global extremum
functions have multiple. Finding a global
extremum is generally difficult; however, many
                                                        local extremum
good methods exist for finding a local
extremum: an extremum for some region
R \subset \mathbb{R}^n.
The method explored here is called gradient
                                                        gradient descent
descent. It will soon become apparent why it
has this name.
Stationary points
Recall from basic calculus that a function f of a
single variable had potential local extrema
where df(x)/dx = 0. The multivariate version of
this, for multivariate function f, is
                     \operatorname{grad} f = 0.
A value X for which Eq. 1 holds is called a
stationary point. However, as in the univariate stationary point
case, a stationary point may not be a local
extremum; in these cases, it called a saddle
                                                         saddle point
Consider the hessian matrix H with values, for hessian matrix
independent variables x_i,
               H_{ij} = \partial_{x_i x_i}^2 f. \tag{2}
For a stationary point X, the second partial
                                                        second partial derivative test
derivative test tells us if it is a local maximum,
local minimum, or saddle point:
 minimum If H(X) is positive definite (all its
                                                        positive definite
      eigenvalues are positive).
       X is a local minimum.
 maximum If H(X) is negative definite (all its
                                                        negative definite
      eigenvalues are negative),
      X is a local maximum.
and negative eigenvalues),
      X is a saddle point.
These are sometimes called tests for concavity:
minima occur where f is convex and maxima
where f is concave (i.e. where -f is convex).
It turns out, however, that solving Eq. 1 directly
for stationary points is generally hard.
Therefore, we will typically use an iterative
technique for estimating them.
The gradient points the way
Although Eq. 1 isn't usually directly useful for
computing stationary points, it suggests
iterative techniques that are. Several techniques
rely on the insight that the gradient points
                                                        the gradient points toward stationary points
toward stationary points. Recall from
Lec. vecs.grad that grad f is a vector field that
points in the direction of greatest increase in f.
Consider starting at some point x_0 and wanting
to move iteratively closer to a stationary point.
So, if one is seeking a maximum of f, then
choose x_1 to be in the direction of grad f. If one
is seeking a minimum of f, then choose x_1 to be
opposite the direction of grad f.
The question becomes: how far \alpha should we go
in (or opposite) the direction of the gradient?
Surely too-small \alpha will require more iteration
and too-large \alpha will lead to poor convergence or
missing minima altogether. This framing of the
problem is called line search. There are a few
common methods for choosing \alpha, called the step \alpha step size \alpha
size, some more computationally efficient than
others.
Two methods for choosing the step size are
described below. Both are framed as
minimization methods, but changing the sign of
the step turns them into maximization methods.
The classical method
                  g_k=\operatorname{grad} f(x_k),
the gradient at the algorithm's current estimate
x_k of the minimum. The classical method of
choosing \alpha is to attempt to solve analytically for
              \alpha_k = \underset{\sim}{\operatorname{argmin}} f(x_k - \alpha g_k).
This solution approximates the function f as one
varies \alpha. It is approximate because as \alpha varies,
so should x. But even with \alpha as the only
variable, Eq. 4 may be difficult or impossible to
solve. However, this is sometimes called the
"optimal" choice for \alpha. Here "optimality" refers
not to practicality but to ideality. This method is
rarely used to solve practical problems.
The algorithm of the classical gradient descent
method can be summarized in the pseudocode
of Algorithm grad.1. It is described further in
                                                        1. Kreyszig, Advanced Engineering Mathematics, § 22.1.
Kreyszig.<sup>1</sup>
Algorithm grad.1 Classical gradient descent
  1: procedure classical_minimizer(f,x_0,T)
        while \delta x > T do \triangleright until threshold T is met
            g_k \leftarrow \operatorname{grad} f(x_k)
            \alpha_k \leftarrow \operatorname{argmin}_{\alpha} f(x_k - \alpha g_k)
            x_{k+1} \leftarrow x_k - \alpha_k g_k
            \delta \mathbf{x} \leftarrow \|\mathbf{x}_{k+1} - \mathbf{x}_k\|
           k \leftarrow k+1
        end while
        return \ x_k \qquad \triangleright \ the \ threshold \ was \ reached
 10: end procedure
The Barzilai and Borwein method
In practice, several non-classical methods are
used for choosing step size \alpha. Most of these
construct criteria for step sizes that are too small
and too large and prescribe choosing some \alpha
that (at least in certain cases) must be in the
sweet-spot in between. Barzilai and Borwein<sup>2</sup>
                                                        2. Jonathan Barzilai and Jonathan M. Borwein. ?Two-Point Step Size
                                                        Gradient Methods? inIMA Journal of Numerical Analysis: 8.1 (january 1988), pages 141–148. issn: 0272-4979. doi: 10.1093/imanum/8.1.
developed such a prescription, which we now
Let \Delta x_k = x_k - x_{k-1} and \Delta g_k = g_k - g_{k-1}. This
                                                                 \alpha_k = \frac{dg}{g} = \frac{d}{g}
method minimizes \|\Delta x - \alpha \Delta g\|^2 by choosing
                  lpha_{
m k} = rac{\Delta x_{
m k} \cdot \Delta g_{
m k}}{\Delta g_{
m k} \cdot \Delta g_{
m k}}.
The algorithm of this gradient descent method
can be summarized in the pseudocode of
Algorithm grad.2. It is described further in
Barzilai and Borwein.<sup>3</sup>
 Algorithm grad.2 Barzilai and Borwein gradient
descent
 1: procedure barzilai_minimizer(f,x<sub>0</sub>,T)
  2: while \delta x > T do \triangleright until threshold T is met
           g_k \leftarrow \operatorname{grad} f(x_k)
           \Delta g_k \leftarrow g_k - g_{k-1}
           \begin{array}{l} \Delta x_k \leftarrow x_k - x_{k-1} \\ \alpha_k \leftarrow \frac{\Delta x_k \cdot \Delta g_k}{\Delta g_k \cdot \Delta g_k} \end{array}
           x_{k+1} \leftarrow x_k - \alpha_k g_k
            \delta \mathbf{x} \leftarrow \|\mathbf{x}_{k+1} - \mathbf{x}_k\|
           k \leftarrow k+1
 10: end while
 11: \operatorname{return} x_k > \operatorname{the threshold was reached}
 12: end procedure
Example opt.grad-1
                                                        re: Barzilai and Borwein gradient descent
Consider the functions (a) f_1 : \mathbb{R}^2 \to \mathbb{R} and (b)
 f_2:\mathbb{R}^2 \to \mathbb{R} defined as
         f_1(\mathbf{x}) = (x_1 - 25)^2 + 13(x_2 + 10)^2 (6)
       f_2(x) = \frac{1}{2}x \cdot Ax - b \cdot x \tag{7}

    Use the method of Barzilai and Borwein<sup>a</sup>

starting at some x_0 to find a minimum of each
 function.
  a. Barzilai and Borwein, ?Two-Point Step Size Gradient Methods?
 First, load some Python packages.
  from sympy import *
  import numpy as np
  import matplotlib.pyplot as plt
  from IPython.display import display, Markdown, Latex
  from tabulate import tabulate
 We begin by writing a class
 gradient_descent_min to perform the
 gradient descent. This is not optimized for
 speed.
  class gradient_descent_min():
     """ A Barzilai and Borwein gradient descent class.
     * f: Python function of x variables
      * x: list of symbolic variables (eg [x1, x2])
      * x0: list of numeric initial guess of a min of
      * T: step size threshold for stopping the descen
     To execute the gradient descent call descend method
     nb: This is only for gradients in cartesian
        coordinates! Further work would be to impleme
        this in multiple or generalized coordinates.
        See the grad method below for implementation
     def __init__(self,f,x,x0,T):
      self.x = Array(x)
      self.x0 = np.array(x0)
      self.T = T
      self.n = len(x0) # size of x
      self.g = lambdify(x,self.grad(f,x),'numpy')
      self.xk = np.array(x0)
self.table = {}
     def descend(self):
      x0 = self.x0
      T = self.T
      g = self.g
      # initialize variables
      x_k = x0
      dx = 2*T # can't be zero
      x_km1 = .9*x0-.1 # can't equal x0
      g_{km1} = np.array(g(*x_km1))
      table_data = [[N,x0,np.array(g(*x0)),0]]
      while (dx > T and N < N_max) or N < 1:
        N += 1 # increment index
        g_k = np.array(g(*x_k))
        dg_k = g_k - g_{m1}
        dx_k = x_k - x_{m1}
        {\tt alpha\_k = abs(dx\_k.dot(dg\_k)/dg\_k.dot(dg\_k))}
        x_k = x_k - alpha_k*g_k
        table_data.append([N,x_k,g_k,alpha_k])
        self.xk = np.vstack((self.xk,x_k))
        g_km1 = g_k
        dx = np.linalg.norm(x_k - x_km1) # check
      self.tabulater(table_data)
     def tabulater(self,table_data):
      np.set_printoptions(precision=2)
      tabulate.LATEX_ESCAPE_RULES={}
      self.table['python'] = tabulate(
        \label{eq:headers=["N","x_k","g_k","alpha"],} headers=["N","x_k","g_k","alpha"],
      self.table['latex'] = tabulate(
        table_data,
        headers=[
          "$N$","$\ \x^x,"$\ \g}_k$","$\\alpha
        tablefmt="latex_raw",
     def grad(self,f,x): # cartesian coord's gradient
      return derive_by_array(f(x),x)
  First, consider f<sub>1</sub>.
  var('x1 x2')
  x = Array([x1,x2])
  f1 = lambda x: (x[0]-25)**2 + 13*(x[1]+10)**2
  gd = gradient_descent_min(f=f1,x=x,x0=[-50,40],T=1e-
 Perform the gradient descent.
 gd.descend()
 Print the interesting variables.
  print(gd.table['python'])
   0 [-50 40]
                         [-150 1300]
   1 [-43.65 -15.03] [-150 1300]
                                                   0.0423296
  2 [-38.36 -10.] [-137.3 -130.74]
                                                   0.0384979
   3 [-33.11 -10.] [-1.27e+02 1.24e-01]
                                                  0.041454
   4 [ 24.99 -10. ] [-1.16e+02 -9.62e-03] 0.499926
   5 [ 25. -10.05]
                         [-0.02 0.12]
                                                    0.499999
   6 [ 25. -10.]
                         [-1.84e-08 -1.38e+00] 0.0385225
  7 [ 25. -10.]
                         [-1.70e-08 2.19e-03] 0.0384615
  8 [ 25. -10.]
                         [-1.57e-08 0.00e+00] 0.0384615
  Now let's lambdify the function f1 so we can
  Now let's plot a contour plot with the gradient
  descent overlaid.
  fig, ax = plt.subplots()
  X1 = np.linspace(-100,100,100)
 X2 = np.linspace(-50,50,100)
X1, X2 = np.meshgrid(X1,X2)
  F1 = f1_lambda(X1,X2)
  plt.contourf(X1,X2,F1)
  plt.colorbar()
   # gradient descent plot
  from mpl_toolkits.mplot3d import Axes3D
  from matplotlib.collections import LineCollection
  xX1 = gd.xk[:,0]
  xX2 = gd.xk[:,1]
  points = np.array([xX1, xX2]).T.reshape(-1, 1, 2)
  segments = np.concatenate(
    [points[:-1], points[1:]], axis=1
   lc = LineCollection(
     segments,
     cmap=plt.get_cmap('Reds')
  lc.set_array(np.linspace(0,1,len(xX1))) # color segs
  lc.set_linewidth(3)
  ax.autoscale(False) # avoid the scatter changing lim
  ax.add_collection(lc)
  ax.scatter(
   xX1,xX2,
   zorder=1,
    marker="o",
    color=plt.cm.Reds(np.linspace(0,1,len(xX1))),
     edgecolor='none'
  Now consider f<sub>2</sub>.
                                                                             Figure grad.1:
  A = Matrix([[10,0],[0,20]])
  b = Matrix([[1,1]])
  def f2(x):
   X = Array([x]).tomatrix().T
   return 1/2*X.dot(A*X) - b.dot(X)
  gd = gradient_descent_min(f=f2,x=x,x0=[50,-40],T=1e-8
  Perform the gradient descent.
  gd.descend()
 Print the interesting variables.
  print(gd.table['python'])
   0 [50-40]
                      [ 499. -801.]
                                                           0
   1 [17.58 12.04] [499. -801.]
                                                 0.0649741
   2 [8.07-1.01] [174.78 239.88]
                                                 0.0544221
   3 [3.62 0.17] [79.66 -21.22]
                                                 0.0558582
   4 [ 0.49 -0.05] [35.16 2.49]
                                                 0.0889491
   5 [0.1 0.14]
                       [ 3.89 -1.94]
                                                 0.0990201
                                                 0.0750849
   6 [0.1 0. ]
                        [0.04 \ 1.9]
   7 [0.1 0.05]
                       [ 0.01 -0.95]
                                                  0.050005
   8 [0.1 0.05]
                       [4.74e-03 9.58e-05]
                                                0.0500012
   9 [0.1 0.05]
                       [ 2.37e-03 -2.38e-09] 0.0999186
  10 [0.1 0.05]
                       [1.93e-06 2.37e-09]
                                                        0.1
                       [ 0.00e+00 -2.37e-09] 0.0999997
  11 [0.1 0.05]
  Now let's lambdify the function f2 so we can
  f2_lambda = lambdify((x1,x2),f2(x),'numpy')
 Now let's plot a contour plot with the gradient
 descent overlaid.
 fig, ax = plt.subplots()
  # contour plot
  X1 = np.linspace(-100,100,100)
  X2 = np.linspace(-50,50,100)
  X1, X2 = np.meshgrid(X1,X2)
  F2 = f2_lambda(X1,X2)
  plt.contourf(X1,X2,F2)
  plt.colorbar()
   # gradient descent plot
  from mpl_toolkits.mplot3d import Axes3D
  from matplotlib.collections import LineCollection
  xX1 = gd.xk[:,0]
  xX2 = gd.xk[:,1]
  points = np.array([xX1, xX2]).T.reshape(-1, 1, 2)
  segments = np.concatenate(
     [points[:-1], points[1:]], axis=1
   lc = LineCollection(
     segments,
     cmap=plt.get_cmap('Reds')
```

lc.set_array(np.linspace(0,1,len(xX1))) # color segs

ax.autoscale(False) # avoid the scatter changing li

color=plt.cm.Reds(np.linspace(0,1,len(xX1))),

Python code in this section was generated from a Jupyter notebook named gradient_descent.ipynb with a python3 Figure grad.2:

lc.set_linewidth(3)

ax.add_collection(lc) ax.scatter(xX1,xX2, zorder=1, marker="o",

edgecolor='none'

plt.show()

opt.grad Gradient descent

Let