

1 ► Accelerated Gradient Descent

Prof. D. Kressner
M. Steinlechner

- a) We consider the overdetermined linear system of equations $A\mathbf{x}^* + \mathbf{w} = \mathbf{b}$ with $A \in \mathbb{R}^{n \times p}$ a matrix with full column rank, $\mathbf{x} \in \mathbb{R}^p$ the unknown solution, and $\mathbf{w} \in \mathbb{R}^n$ denotes unknown noise. The least squares estimator for the unknown solution is then given by

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} f(\mathbf{x}), \quad \text{with } f(\mathbf{x}) := \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

Show that f is a μ -strongly convex function with μ given by the smallest eigenvalue of $A^T A$. Furthermore, show that the gradient of f is Lipschitz continuous with the smallest Lipschitz constant L given by the largest eigenvalue of $A^T A$.

We obtain the Hessian $H(x) = A^T A$, hence f is μ -strongly convex with

$$\mu = \lambda_{\min}(A^T A),$$

and the gradient of f is Lipschitz continuous with Lipschitz constant

$$L = \lambda_{\max}(A^T A).$$

- b) Implement Algorithm 4.28, *Accelerated Gradient Descent*. Use the provided template on the lecture homepage, `AGD.m`

```

1  function [X,fX] = AGD( f, df, x0, mu, L, opts );
   % load default values for option parameters if not set by user
3  if ~exist( 'opts', 'var' );    opts = struct();    end
   if ~isfield( opts, 'tol' );   opts.tol = 1e-6;    end
5  if ~isfield( opts, 'maxiter' ); opts.maxiter = 1000; end

7  X(:,1) = x0;
   fX(:,1) = f(x0);
9
   xk = x0;
11  yk = x0;
   gk = df(x0);
13
   kappa = sqrt(mu / L);
15  kappa = (1 - kappa) / (1 + kappa);

17  k = 1;
   while norm(gk) > opts.tol && k < opts.maxiter
19     % perform step
       yk_new = xk - gk / L ;
21     xk = yk_new + kappa * (yk_new - yk);
       yk = yk_new;
23
       % prepare next step
25     gk = df(xk);
       X(:,k+1) = xk;
27     fX(:,k+1) = f(xk);

29     k = k + 1;
   end
31 end
end
```

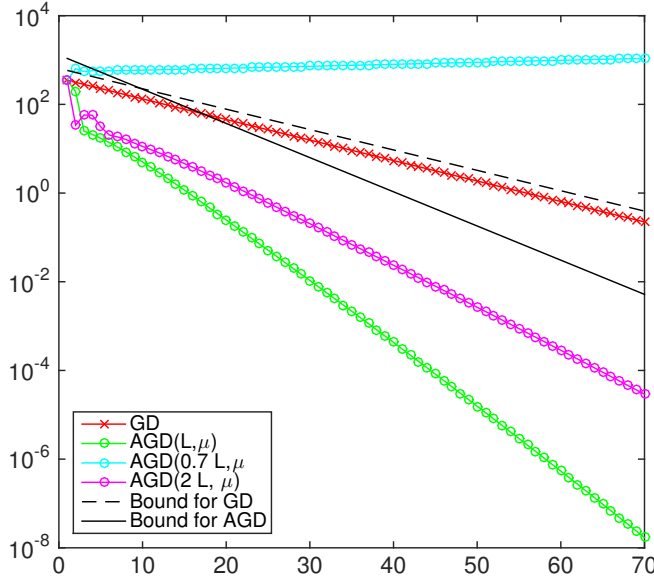
- c) We will now investigate the convergence rates of standard Gradient Descent (GD) and Accelerated Gradient Descent (AGD) when applied to the cost function f . An implementation of GD is available on the lecture homepage, `GD.m`. Using

the provided template `ex9problem1.m`, plot the convergence rates of GD for the optimal choice $\alpha = \frac{2}{L+\mu}$ which exploits the strong convexity of the function. Compare to the convergence rate of AGD using your implementation from **b)**. Try out how the convergence speed changes as you use non-optimal values for L and μ .

```

2   rng(11);
3   % set up least squares problem
4   n = 500;
5   p = 10;
6   A = rand( n, p );
7   x = rand( p, 1 );
8   noise = 1e-5*randn( n, 1 );
9   b = A*x + noise;
10
11  % obtain the exact minimizer by direct least squares solver (as the noise is gaussian)
12  % Note that due to the noise, x is not the same as x_exact
13  x_exact = A\b;
14
15  % Definition of function, derivative and Hessian
16  f = @(x) 0.5 * norm( A*x - b )^2;
17  df = @(x) A' * (A*x - b);
18  H = A'*A;
19
20  fx_exact = f(x_exact);
21
22  ev = eig(H);
23
24  % Compute Lipschitz constant L and strong convexity parameter mu exactly
25  L = max(ev)
26  mu = min(ev);
27
28  opts.maxiter = 70;
29  opts.tol = 1e-6;
30
31  x0 = rand( p, 1 );
32
33  [X_GD, fX_GD] = GD( f, df, x0, 2/(L+mu), opts );
34  [X_AGD, fX_AGD] = AGD( f, df, x0, mu, L, opts );
35  % Run AGD again with "wrong" constants
36  [X_AGD_07, fX_AGD_07] = AGD( f, df, x0, mu, 0.7*L, opts );
37  [X_AGD_2, fX_AGD_2] = AGD( f, df, x0, mu, 2*L, opts );
38
39  % Check for convergence:
40  norm( x_exact - X_GD(:,end) )
41  norm( x_exact - X_AGD(:,end) )
42  norm( x_exact - X_AGD_07(:,end) )
43  norm( x_exact - X_AGD_2(:,end) )
44
45  semilogy( 1:length(fX_GD), fX_GD - fx_exact, '-xr' )
46  hold on
47  semilogy( 1:length(fX_AGD), fX_AGD - fx_exact, '-og' )
48  semilogy( 1:length(fX_AGD_07), fX_AGD_07 - fx_exact, '-oc' )
49  semilogy( 1:length(fX_AGD_2), fX_AGD_2 - fx_exact, '-om' )
50
51  set(gca,'fontsize',16)
52
53  k = 1:opts.maxiter;
54  kappa = L / mu;
55  gamma = (kappa-1) / (kappa+1);
56  semilogy( k, 0.5*L*gamma.^(2*k), '--k' )
57  semilogy( k, L*(1-sqrt(mu/L)).^k, '-k' )
58
59  l = legend( 'GD', 'AGD(L,\mu)', 'AGD(0.7_L,\mu)', 'AGD(2_L,\mu)', 'Bound_for_GD', 'Bound_for_AGD' )
60  set(l,'location','SouthWest')
61
62

```



2 ► Binary logistic regression, part 2

In the last exercise sheet, we have seen the *binary logistic regression* approach. Let $\mathbf{a}_i \in \mathbb{R}^n$ be sampling points and b_i be the associated binary class labels. Then, we want to determine the maximum log-likelihood estimator

$$\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}), \quad \text{with } f(\mathbf{x}) = - \sum_{i=1}^n \log(h(b_i \mathbf{a}_i^\top \mathbf{x})),$$

where $h(t) = 1/(1 + \exp(t))$ is the sigmoid function. In the last exercise sheet, we have shown that the function is convex, but in general not strongly convex, depending on the data matrix $A = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_n]^\top$.

Hence, we introduce a regularization term, such that the minimization problem is now given as

$$x^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f_\sigma(\mathbf{x}), \quad \text{with } f_\sigma(\mathbf{x}) = f(\mathbf{x}) + \frac{\sigma}{2} \|\mathbf{x}\|_2^2.$$

a) Show that the regularized function f_σ is σ -strongly convex.

The Hessian of the regularized function f_σ is given by

$$H_\sigma = H(x) + \sigma I,$$

where $H(x)$ is the Hessian of the non-regularized function f with which we worked in the last exercise. There, it was shown to be positive semi-definite. Hence, if $H(x)$ has eigenvalues λ_i , then H_σ has the shifted eigenvalues $\lambda_i + \sigma$. Therefore, the f_σ is μ -strongly convex with parameter $\mu = \sigma$.

b) Show that the gradient of f_σ is Lipschitz continuous with $L = \frac{1}{4} \|A\|_2 + \sigma$.

With the same argument as in a), the maximum eigenvalue of $H_\sigma(x)$ is the maximum eigenvalue of the non-regularized Hessian $H(x)$ shifted by the regularization parameter σ .

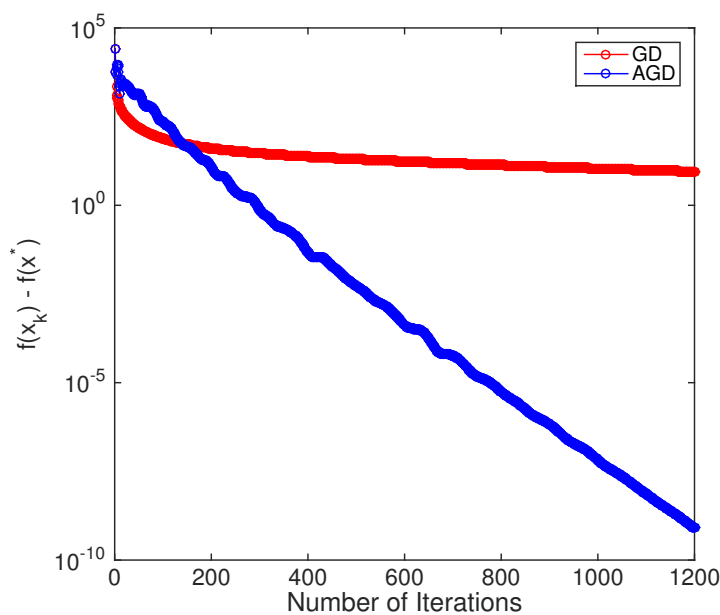
c) We will now test the performance of GD and AGD applied to binary logistic regression. For simplicity, we choose the regularization parameter $\sigma = 1$. As a dataset, we will use the **a4a** dataset from <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>, a processed version of the “Adult” dataset. This dataset aims to predict from 4781 samples from the 1994 US census the probability of a person to earn over 50000 USD per year given 123 attributes

such as age, sex, education, etc. Complete the template `ex9problem2.m` using your implementation of AGD from the previous question and the prepared dataset `adult.mat` from the lecture homepage.

```

2  % Load the dataset to get matrices A, b and the exact
   % solution xexact
   Load('adult.mat')
4
6  % Regularization parameter
   sigma = 1;
8
8  % Lipschitz constant
   L = 0.25*svds(A,1)^2 + sigma;
10
12 % strong convexity parameter
   mu = sigma;
14
14 h = @(t) 1./(1 + exp(-t));
   f = @(x) -sum(log(h(b.*(A*x)))) + 0.5*sigma*x'*x;
16 df = @(x) -A' * ( (1-h(b.*(A*x))) .* b ) + sigma*x;
18
18 fxexact = f(xexact)
20
20 opts.maxiter = 1200;
   opts.tol = 1e-6;
22
22 x0 = rand( size(A,2), 1 );
24
24 [X_GD, fX_GD ] = GD ( f, df, x0, 2/(L+mu) , opts );
26 [X_AGD, fX_AGD] = AGD( f, df, x0, mu, L , opts );
28
28 semilogy( 1:length(fX_GD), fX_GD - fxexact, '-or' )
30 hold on
   set(gca,'fontsize',16)
32 semilogy( 1:length(fX_AGD), fX_AGD - fxexact, '-ob' )
34 xlabel('Number_of_Iterations')
   ylabel('f(x_k) - f(x^*)')
36 legend('GD','AGD')

```



3 ► Tangent cones

Follows directly from the definition. See also discussion in the exercise class.