

OLS with ℓ_1 and ℓ_2 regularization

CEE 629. System Identification

Duke University, Fall 2017

ℓ_1 regularization

- The ℓ_1 norm of a vector $v \in \mathbb{R}^n$ is given by $\|v\|_1 = \sum |v_i|$
The gradient of $\|v\|_1$ is not defined if an element of v is zero.
- In ℓ_1 regularization, the objective $J(a) = \|y - f(y; a)\|_2^2$ is penalized with a term $\alpha \|a\|_1$, where α is called the *regularization parameter*, or the *penalty factor*.
- The effect of ℓ_1 regularization is to force some of the model parameters, a_i , to zero (exactly). If the parameters are coefficients for bases of the model, then ℓ_1 regularization is a means to remove un-important bases of the model. It is a form of model reduction
- It is applicable to linear and nonlinear models.
- It has similar but not identical effects as ℓ_2 regularization and truncated SVD solutions
- a.k.a. LASSO, Compressed Sensing, Support Vector Machine

But ℓ_1 is not smooth at the optimum point.

What can one say about necessary conditions for optimality?

ℓ_1 regularization can be recast as a Quadratic Program (QP)

$$J(a) = \|y - Xa\|_2^2 + \alpha \|a\|_1 = \|y - Xa\|_2^2 + \alpha \sum_i |a_i| \quad (1)$$

\Leftrightarrow

$$a_i \equiv p_i - q_i \quad (2)$$

$$|a_i| \equiv p_i + q_i \quad (3)$$

$$J(p, q) = \|y - Xp + Xq\|_2^2 + \alpha \sum_i (p_i + q_i) \quad (4)$$

such that $p_i \geq 0$ and $q_i \geq 0 \quad \forall i$

Five conditions to show that ℓ_1 is a QP

1. $p_i - q_i = a_i$
2. $p_i + q_i = |a_i|$
3. $p_i \geq 0 \quad q_i \geq 0$
 - (1)&(2) \Rightarrow (3)&(4)
4. $a_i > 0 \Rightarrow p_i = a_i, \quad q_i = 0$
 - (1)&(3) \nRightarrow (2)&(4)
- $a_i < 0 \Rightarrow p_i = 0, \quad q_i = -a_i$
 - (1)&(3) \Rightarrow (4)&(5)
- $a_i = 0 \Rightarrow p_i = 0, \quad q_i = 0$
5. $\min(p_i + q_i) = a_i$

Solving the ℓ_1 regularization QP

$$J(p, q) = \|y - Xp + Xq\|_2^2 + \alpha \sum_i (p_i + q_i)$$

such that $p_i \geq 0$ and $q_i \geq 0 \quad \forall i$

But which elements of p_i and which elements of q_i should be set to zero?
... (the “active” constraints)

This question does not have an answer in closed form.

So we solve the problem incrementally.

Incremental formulation

$$p^{(k+1)} = p^{(k)} + u^{(k)}, \quad p^{(k)} + u^{(k)} \geq 0 \tag{5}$$

$$q^{(k+1)} = q^{(k)} + v^{(k)}, \quad q^{(k)} + v^{(k)} \geq 0 \tag{6}$$

Lagrangian

$$L(p, q, \mu, \nu) = J(p, q) + \mu_{i \in A_u}^{(k)\top} (p^{(k)} + u^{(k)})_{i \in A_u} + \nu_{i \in A_v}^{(k)\top} (q^{(k)} + v^{(k)})_{i \in A_v}$$

- k is the iteration number
- A_u is the set of active equality constraints $p_i = 0$
- A_v is the set of active equality constraints $q_i = 0$
- $\mu_{i \in A_u}^{(k)}$ is the set of Lagrange multipliers for the active constraint set A_u
- $\nu_{i \in A_v}^{(k)}$ is the set of Lagrange multipliers for the active constraint set A_v

Minimize the Lagrangian w.r.t. (p, q) ,

Maximize the Lagrangian w.r.t. (μ, ν)

$$\begin{aligned}
 L = & y^\top y - 2(p + u)^\top X^\top y + p^\top \alpha \\
 & + 2(q + v)^\top X^\top y + q^\top \alpha \\
 & + (p + u)^\top X^\top X (p + u) \\
 & - 2(q + v)^\top X^\top X (p + u) \\
 & + \mu_{i \in A_u}^\top (p + u)_{i \in A_u} \\
 & + \nu_{i \in A_v}^\top (q + v)_{i \in A_v}
 \end{aligned}$$

The necessary conditions for optimality ...

$$\frac{\partial L}{\partial u} = 0 : -2X^\top y + \alpha + \mu_{i \in A_u} + 2X^\top X p + 2X^\top X u - 2X^\top X q - 2X^\top X v = 0$$

$$\frac{\partial L}{\partial v} = 0 : +2X^\top y + \alpha + \nu_{i \in A_v} - 2X^\top X p - 2X^\top X u + 2X^\top X q + 2X^\top X v = 0$$

$$\frac{\partial L}{\partial \mu_{i \in A_u}} = 0 : (p + u)_{i \in A_u} = 0$$

$$\frac{\partial L}{\partial \nu_{i \in A_v}} = 0 : (q + v)_{i \in A_v} = 0$$

... are linear in our unknowns u, v, μ, ν ...

$$\begin{bmatrix}
 2X^\top X & -2X^\top X & I_{i \in A_u}^\top & & \\
 -2X^\top X & 2X^\top X & & I_{i \in A_v}^\top & \\
 & I_{i \in A_u} & & & \\
 & & I_{i \in A_v} & &
 \end{bmatrix}
 \begin{bmatrix}
 u \\
 v \\
 \mu \\
 \nu
 \end{bmatrix}
 =
 \begin{bmatrix}
 2X^\top y - 2X^\top X p + 2X^\top X q - \alpha \\
 -2X^\top y + 2X^\top X p - 2X^\top X q - \alpha \\
 -p_{i \in A_u} \\
 -q_{i \in A_v}
 \end{bmatrix}$$

The number of columns of the matrix $I_{i \in A_u}$ is the number of model parameters. The number of rows of $I_{i \in A_u}$ is the number of active constraints on p . If index i is the element j in the set of active constraints, A_u , then the i, j element of $I_{i \in A_u}$ equals 1.

A trick to block an incremental step from going negative.

Define the index j s.t. $p_j + u_j = \min(p + u) \dots (p_j \geq 0)$

If $p_j + u_j < 0$, then define $\delta_u = -p_j/u_j$ ($\delta_u > 0$)

and adjust the increment by the factor δ_u

$$p^{(k+1)} = p^{(k)} + \delta_u u^{(k)}$$

Likewise ...

Define the index j s.t. $q_j + v_j = \min(q + v) \dots (q_j \geq 0)$

If $q_j + v_j < 0$, then define $\delta_v = -q_j/v_j$ ($\delta_v > 0$)

and adjust the increment by the factor δ_v

$$q^{(k+1)} = q^{(k)} + \delta_v v^{(k)}$$

Implementation in `L1_fit.m`

Example 1, 2-term power polynomial

$$\hat{y}(x; a_1, a_2) = a_1x + a_2x^2 \quad y = \hat{y} + \eta \quad \sigma_\eta = 0.5$$

$$(a_1 = 4, a_2 = 8)$$

OLS criterion

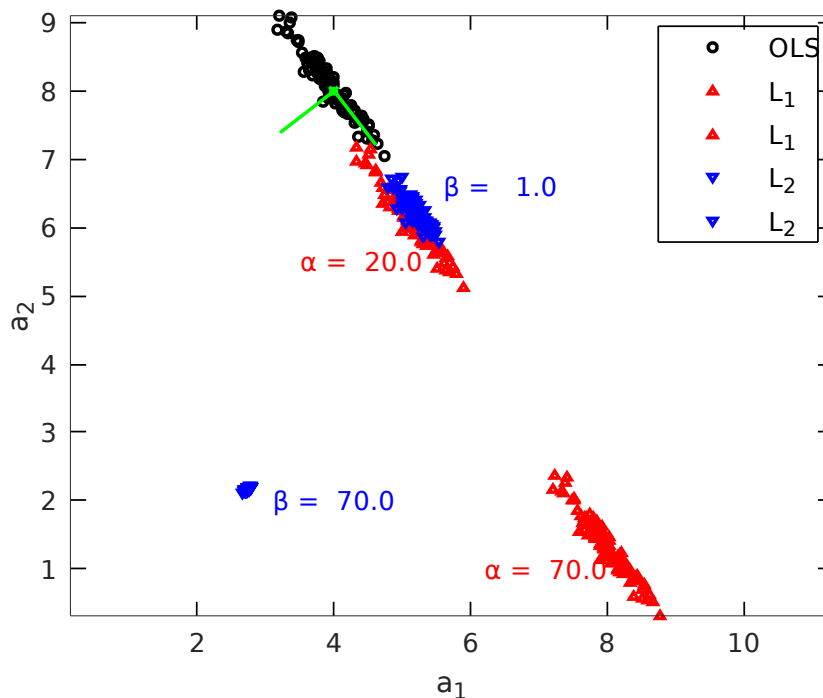
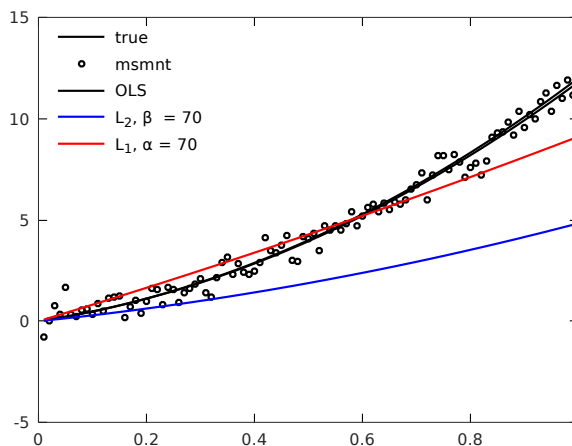
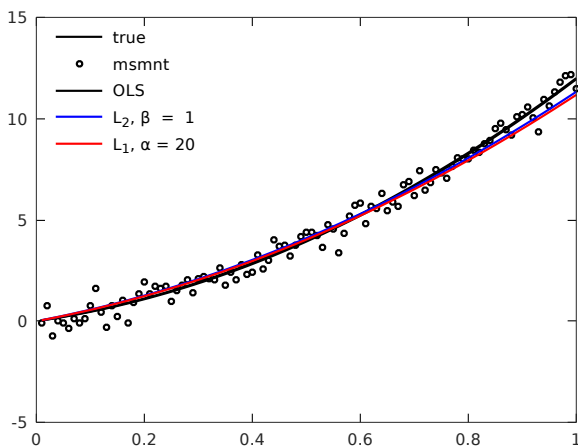
$$\min_a \|y - Xa\|_2^2$$

l_1 regularization criterion

$$\min_a \|y - Xa\|_2^2 + \alpha \|a\|_1$$

l_2 regularization criterion

$$\min_a \|y - Xa\|_2^2 + \beta \|a\|_2$$



Example 2, 100-term Prony Series

A Prony series is a model for the behavior of a dynamic system expressed in terms of a convolution of the system forcing with a kernel, in which the kernel is a series of exponentials.

$$y(t) = \int_0^t h(r)u(t-r) dr$$

where, for example, the kernel can be parameterized as

$$h(r) = k_0 + \sum_{i=1}^n k_i \exp[-r/\tau_i]$$

Prony series are linear in the coefficients k_i but are highly nonlinear in the time constants τ_i .

In the frequency domain, the complex modulus is obtained from the Fourier transform of the convolution kernel,

$$\hat{H}(\omega) = k_0 + \sum_{i=1}^n k_i \frac{i\omega\tau_i}{1 + i\omega\tau_i}$$

The real part of $\hat{H}(\omega)$ is called the *storage modulus* and represents the energy-conserving elasticity of the system. The imaginary part of $\hat{H}(\omega)$ is called the *loss modulus* and relates to the energy dissipation aspects. For linear viscoelastic materials the loss modulus approaches zero as ω approaches zero.

Using ℓ_1 regularization, a Prony series with a very large set of time constants can be fit to complex-modulus data. By tuning the ℓ_1 penalty factor α to be large enough, a large number of coefficients in the Prony series will be forced to zero. Since negative values of the coefficients are not permissible, this problem is a perfect candidate for ℓ_1 regularization.

The system model is then given by

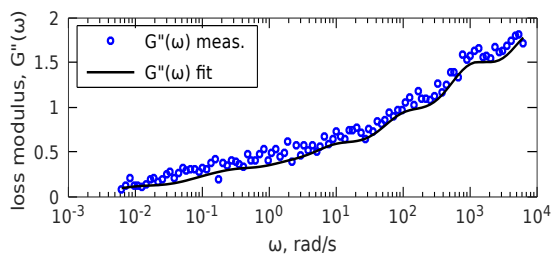
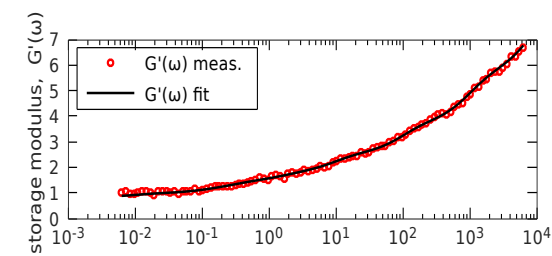
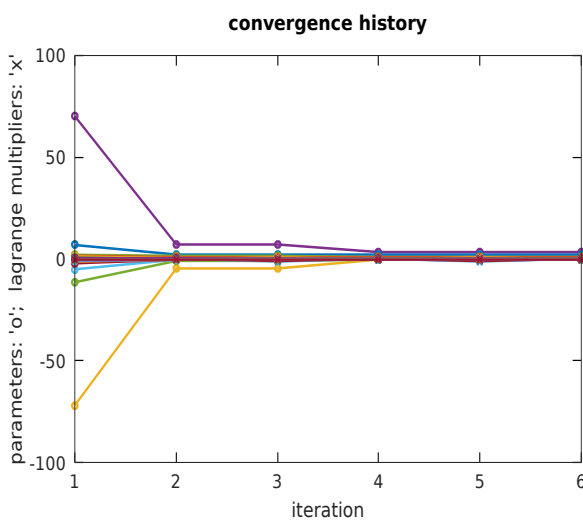
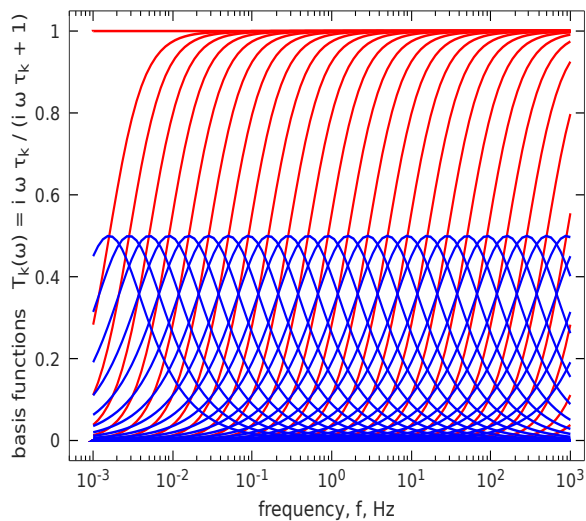
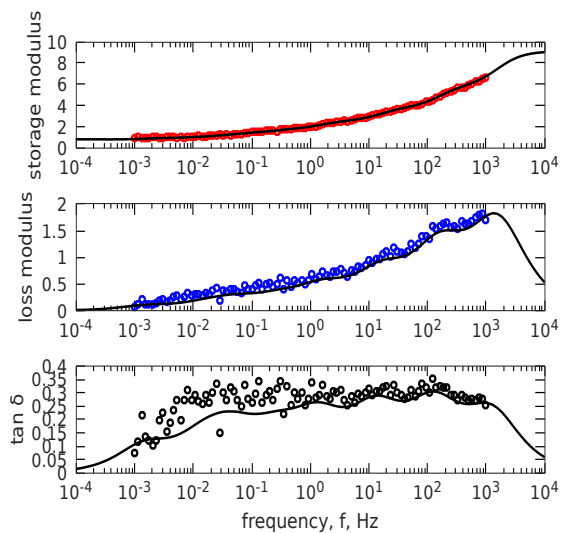
$$\hat{H}(\omega; k) = T(\tau)k$$

where the i -th column of the matrix $T(\tau)$ is $i\omega\tau_i/(1 + i\omega\tau_i)$ and the vector k contains the $n + 1$ coefficients k_0, \dots, k_n .

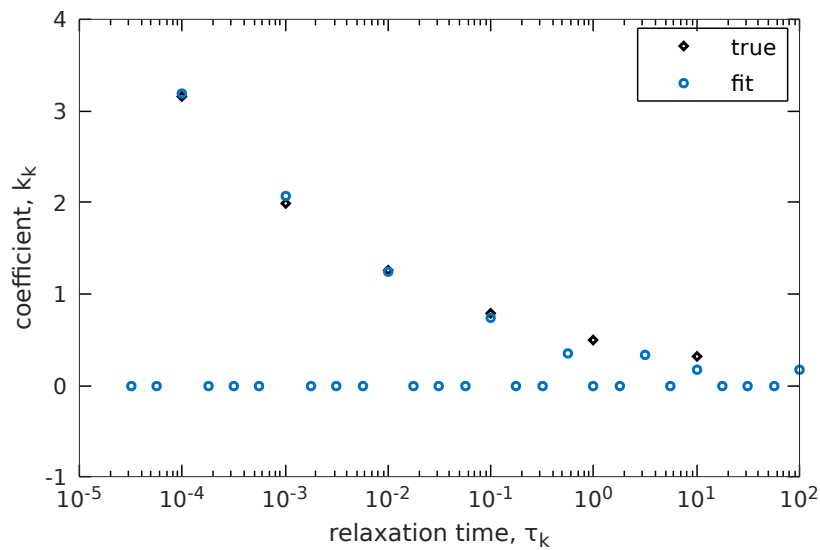
The regularized objective function is

$$J(k) = \|H(\omega) - T(\tau)k\|_2^2 + \alpha \sum_i k \quad \text{s.t. } k \geq 0$$

This is a special (more simple) case of the ℓ_1 objective analyzed in the previous section. The solution is implemented in `PronyFit.m`.



relaxation spectrum



L1_reg_example.m

```

1 % L1_reg_example.m
2 % HPG, Duke Univ. 2013-09, 2017-09
3
4 % generate some data (x,y), noisy in y
5
6 m = 100; % number of measurement points
7 x = [ 1:m]'/m; % independent variables – noise-free
8
9 a = [ 4 8 ]'; % "true" parameter coefficients
10 n = length(a);
11
12 X = zeros(m,n);
13 for i=1:n
14     X(:,i) = x.^i;
15 end
16
17 [eigvec,eigval] = eig(X'*X); % principal directions of parameter covariance
18
19 y = X*a;
20
21 alpha_set = [ 20.0 70.0 ]; % for L1 regularization
22 beta_set = [ 1.0 70.0 ]; % for L2 regularization
23
24 N_MCS = 100;
25 N_R = length(alpha_set);
26
27 a_hat_ols_rec = zeros(n,N_MCS);
28 a_hat_alpha_rec = zeros(n,N_R,N_MCS);
29 a_hat_beta_rec = zeros(n,N_R,N_MCS);
30
31
32 for j=1:N_MCS; % ... start simulation loop
33
34     y_dat = y + 0.5*randn(m,1); % add noise to true data
35
36     a_hat_ols = X \ y_dat; % ordinary least squares
37     y_hat_ols = X*a_hat_ols;
38
39     a_hat_ols_rec(:,j) = a_hat_ols;
40
41     for k=1:N_R
42
43         alpha = alpha_set(k);
44         beta = beta_set(k);
45
46         a_hat_alpha = L1_fit ( X, y_dat, alpha ); % L1 regularization
47         y_hat_alpha = X * a_hat_alpha;
48         a_hat_alpha_rec(:,j,k) = a_hat_alpha;
49
50         a_hat_beta = [ X'*X + beta*eye(n) ] \ X'*y_dat; % L2 regularization
51         y_hat_beta = X*a_hat_beta;
52         a_hat_beta_rec(:,j,k) = a_hat_beta;
53
54     figure(1)
55     clf
56     plot( x,y,'-k', x,y_dat,'ok', ...
57           x,y_hat_ols,'-k', x,y_hat_beta,'-b', x,y_hat_alpha,'-r' )
58     legend('true','msmnt','OLS', ...
59           sprintf('L2, \\beta = %2.0f',beta), ...
60           sprintf('L1, \\alpha = %2.0f',alpha),'location','northwest');
61     legend('boxoff')
62     drawnow
63
64     end
65
66 end % ... end simulation loop

```


L1_fit.m

```

1  function [a, mu, nu, cvg_hst] = L1_fit ( X, y, alpha )
2  % [a, mu, nu, cvg_hst] = L1_fit(X,y,alpha)
3  %
4  % fit model parameters, a, in the model  $\hat{y} = X*a$  to data, y, with
5  % L1 regularization of the parameters.
6  %  $J = || y - X*a ||_{-2}^2 + \alpha * ||a||_{-1}$ 
7  %  $\Leftrightarrow$ 
8  %  $J = || y - X(p-q) ||_{-2}^2 + \alpha * \text{sum}(p+q)$  such that  $p \geq 0, q \geq 0$ 
9  %  $(a = p-q; |a| = p+q)$ 
10 %
11 % INPUT: X ... the design matrix (basis of the model), m x n
12 %        y ... the vector of data to be fit to the model, m x 1
13 %        alpha ... l1 regularization factor for Prony series coefficients
14 %
15 % OUTPUT: a ... the model coefficients
16 %         mu,nu ... Lagrange multipliers for p and q
17 %         cvg_hst ... convergence history
18 %
19
20 %}
21 The main idea behind casting L1 as a QP is that the parameter
22 vector {a} is replaced by the difference of two vectors {a}={p}-{q}
23 that are constrained to be non-negative:  $p_i \geq 0$  for all i and  $q_i \geq 0$ 
24 for all i.
25 If  $a_i > 0$ , then  $p_i = a_i$  and  $q_i = 0$ ;
26 If  $a_i < 0$ , then  $p_i = 0$  and  $q_i = -a_i$ .
27 With this constrained re-parameterization,  $|a_i| = p_i + q_i$ .
28 Note that the dimension of the parameter space doubles, but the KKT equations
29 for the QP are simple and have analytical Hessians and gradients.
30 %}
31
32
33 [m,n] = size(X);
34
35 t = [1:m];
36
37 XtX = 2*X'*X; Xty = 2*X'*y;
38
39 % _good_ initial guess for p and q from non-regularized linear least squares
40 a = XtX \ Xty;
41 p = zeros(n,1); p (find(a > 2*eps)) = a(find(a > 2*eps));
42 q = zeros(n,1); q (find(a < -2*eps)) = -a(find(a < -2*eps));
43
44 MaxIter = 19;
45
46 cvg_hst = zeros ( 5*n, MaxIter ); % convergence history
47
48 for iter = 1:MaxIter
49
50     Au = find(p <= 2*eps); lp = length(Au); % active set for update u
51     Av = find(q <= 2*eps); lq = length(Av); % active set for update v
52
53     Ip = zeros(lp,n); for i=1:lp, Ip(i,Au(i)) = 1; end % constraint gradient u
54     Iq = zeros(lq,n); for i=1:lq, Iq(i,Av(i)) = 1; end % constraint gradient v
55
56     mu = zeros(n,1); nu = zeros(n,1);
57     % KKT equations
58     ITX = [ XtX      -XtX      Ip'      zeros(n,lq) ; % dL/du = 0
59            -XtX      XtX      zeros(n,lp)  Iq'      ; % dL/dv = 0
60            Ip      zeros(lp,n+lp+lq)      ; % dL/d mu = 0
61            zeros(lq,n)  Iq      zeros(lq,lp+lq)  ]; % dL/d nu = 0
62
63     % right-hand-side vector
64     XTY = [ Xty - XtX*p + XtX*q - alpha ; % dL/du = 0
65            -Xty + XtX*p - XtX*q - alpha ; % dL/dv = 0
66            -p(Au) ; % dL/d mu = 0
67            -q(Av) ] ; % dL/d nu = 0
68
69     u_v_mu_nu = ITX \ XTY; % solve the system

```

```

70
71 %
72
73 u      = u_v_mu_nu(1:n);           % update for p
74 v      = u_v_mu_nu(n+1:2*n);      % update for q
75 mu(Au) = u_v_mu_nu(2*n+1:2*n+1p); % what to use mu for?
76 nu(Av) = u_v_mu_nu(2*n+1p+1:2*n+1p+1q); % what to use nu for?
77
78 % if an element of p+u becomes negative, reduce the length of step u
79 du = 1;
80 [p_min, j] = min(p+u);
81 if p_min < 0
82     du = -p(j)/u(j);
83 end
84
85 % if an element of q+v becomes negative, reduce the length of step v
86 dv = 1;
87 [q_min, j] = min(q+v);
88 if q_min < 0
89     dv = -q(j)/v(j);
90 end
91
92 p = p + du*u;
93 q = q + dv*v;
94 a = p - q;
95
96 cvg_hst(:,iter) = [ a ; p ; q ; mu ; nu ]; % convergence history
97
98 % convergence check
99 if ( norm(u) <= norm(p)/1e3 && min(p) > -5*eps && ...
100     norm(v) <= norm(q)/1e3 && min(q) > -5*eps )
101     break;
102 end
103
104 % figure(101); clf; plot(t,y,'o', t,X*a); drawnow;
105
106 end
107
108 cvg_hst = cvg_hst(:,1:iter);
109
110 % figure(102); clf; plot([1:iter], cvg_hst(1:n,:)' ,'-o'); drawnow; %pause(1);
111
112 % _____ L1_fit
113 % H.P. Gavin, 2013-10-04

```

PronyFit_test.m

```

1  % — PronyFit_test.m
2  % test prony series fitting with L1 regularization
3
4  clear all
5
6  % simulated noisy data
7
8  % a relatively small number of 'true' time lags
9  tau = 1e-2*[ 0.01 0.1 1 10 100 1000 ]';
10 ko = 1; % static stiffness
11 k = 0.5 ./ tau.^0.2; % 'true' Prony series coefficients
12
13 points = 100;
14 f_dat = logspace(-3,3,points)';
15 iw = 2*pi*f_dat * sqrt(-1.0);
16
17 % generate complex modulus data with some additive noise and friction effects
18 T = [ ones(points,1) , iw*tau' ./ (iw*tau' + 1.0) ];
19 G_dat = T * [ ko ; k ] + 5e-2*(randn(points,1) + i*randn(points,1));
20
21 friction = 0.1; % definitely NOT linear viscoelasticity!
22 G_dat = real(G_dat) + sqrt(-1)*(friction + imag(G_dat));
23
24 % plot the simulated data
25 figure(1)
26 clf
27 semilogx(f,real(G),'or',f,imag(G),'ob')
28 subplot(311)
29 semilogx(f_dat,real(G_dat),'or')
30 ylabel('storage modulus')
31 subplot(312)
32 semilogx(f_dat,imag(G_dat),'ob')
33 ylabel('loss modulus')
34 subplot(313)
35 semilogx(f_dat,imag(G_dat)./real(G_dat),'ok')
36 ylabel('tan \delta')
37 xlabel('frequency, f_{dat}, Hz')
38
39 % number of terms in the Prony series
40 nID = 97;
41
42 tau_id = logspace(-4.5,2,nID)'; % specified set of time lags, tau
43
44 alpha = 0.50; % l1 regularization parameter
45
46 [ko_id,k_id,cvg_hst] = PronyFit(G_dat,f_dat,tau_id,alpha); % do it!
47
48 f_fit = logspace(-4,4,100);
49 iw_fit = 2*pi*f_fit * sqrt(-1.0);
50 G_fit = ko_id + sum( (k_id.*tau_id * iw_fit) ./ (tau_id * iw_fit + 1) );

```

PronyFit.m

```

1 function [ko,k,cvg_hst] = PronyFit(G_dat,f_dat,tau,alpha)
2 % [ko,k,cvg_hst] = prony_fit(G_dat,f_dat,tau,alpha)
3 %
4 % Fit a prony series to frequency domain complex modulus data
5 %
6 % INPUT:  G_dat ... complex modulus data           Mx1 complex vector
7 %         f_dat ... frequencies where the measured frequency response data is
8 %           evaluated (Hz), (f>0)                 Mx1 real vector
9 %         tau ... specified set of relaxation times, Nx1 real vector
10 %        alpha ... l-1 regularization factor for Prony series coefficients
11 %
12 % OUTPUT: ko, k ... Prony series coefficients
13 %         cvg_hst ... convergence history
14 %
15
16 i = sqrt(-1.0);
17
18 m = length(f_dat);           % number of data points
19 w = 2*pi * f_dat;           % one-sided frequency
20
21 T = [ ones(m,1) , i*w*tau' ./ (i*w*tau' + 1.0) ]; % design matrix
22
23 % plot the basis functions (columns of T) with single-sided log-scale frequencies
24 figure(101)
25 semilogx(f_dat, real(T), '-r', f_dat, imag(T),'-b')
26 ylabel('basis functions T_k(\omega) = i \omega \tau_k / (i \omega \tau_k + 1)')
27 xlabel('frequency, f, Hz')
28 axis( [ 0.5*min(f_dat) 1.5*max(f_dat) -0.05 1.05 ] )
29
30 [m,n] = size(T);           % m = number of data points; n = number of parameters
31
32 TtT = 2*real(T'*T);       % taking real part is like adding complex conjugate
33 TtG = 2*real(T'*G_dat);
34
35 k = TtT \ TtG; % O.L.S. fit is a better initial guess, even though infeasible, k<0
36
37 MaxIter = 20;             % usually enough
38 cvg_hst = zeros(2*n,MaxIter); % convergence of coefficients and multipliers
39
40 for iter = 1:MaxIter      % — Start Main Loop
41
42     A = find(k < 2*eps); l = length(A); % active set for update constraints
43     % Ia is the constraint gradient w.r.t. the step vector h
44     Ia = zeros(l,n); for i=1:l, Ia(i,A(i)) = 1; end
45
46     % KKT equations
47     lambda = zeros(n,1); % Lagrange multiplier
48     XTX = [ TtT Ia' ; % dL / dh
49            Ia zeros(1,1) ]; % dL / d lambda
50
51     XTY = [ -TtT*k + TtG - alpha ; % dL / dh
52            -k(A) ]; % dL / d lambda
53
54     h_lambda = XTX \ XTY; % solve the system
55
56 %
57
58     h = h_lambda(1:n); % the step
59     lambda(A) = h_lambda(n+1:end); % the non-zero multipliers
60
61     % if an element of k+h becomes negative, reduce the length of step h to dh*h
62     dh = 1;
63     [h_test_min, idx] = min(k+h);
64     if h_test_min < 0
65         dh = -k(idx)/h(idx);
66     end
67

```

```

68     k = (k + dh * h);                                     % update the Prony coefficients
69
70     figure(102)                                         % plot msmnt points and fit
71     G_hat = T*k;
72     clf
73     subplot(211)
74     semilogx ( w,real(G_dat),'or', w,real(G_hat),'-k')
75     ylabel('storage modulus, G''(\omega) ')
76     legend('G''(\omega) meas.','G''(\omega) fit','location','northwest')
77     subplot(212)
78     semilogx ( w,imag(G_dat),'ob', w,imag(G_hat),'-k')
79     xlabel('\omega, rad/s')
80     ylabel('loss modulus, G''(\omega) ')
81     legend('G''(\omega) meas.','G''(\omega) fit','location','northwest')
82     drawnow
83
84     cvg_hst(:,iter) = [ k ; lambda ];                   % convergence history
85     if ( norm(h) < norm(k)/1e3 && min(k) > -5*eps )    % convergence check
86         break;
87     end
88
89     end                                               % — End Main Loop
90
91     cvg_hst = cvg_hst(:,1:iter);
92
93     % asymptotic standard errors of the Prony coefficients ???
94     % k_std_err = sqrt(norm(G_dat-T*k)/(m-n+1) * diag(real(inv(TtT))));
95
96     ko = k(1);
97     k = k(2:n);
98
99     % _____ PRONY_FIT
100    % HP Gavin, 2013-10-04

```