

## Math Foundations of ML, Fall 2022

### Homework #5

Due Monday October 31, at 5:00pm ET

As stated in the syllabus, unauthorized use of previous semester course materials is strictly prohibited in this course.

1. In this problem, we will solve a stylized regression problem using the data set `hw05p1_clusterdata.mat`. This file contains (noisy) samples of a function  $f(t)$  for  $t \in [0, 1]$ . In fact, the function is the same as the one we used in Homework 4:

$$f_{\text{true}}(t) = \frac{\sin(12(t + 0.2))}{t + 0.2}.$$

The sample locations are in the vector `T`; the sample values are in `y`. If you plot these, you will see that the samples come in four clusters. We are going to use nonlinear regression using the orthobasis of Legendre polynomials — in the notes, these polynomials were defined on  $[-1, 1]$ , we will use the analogous functions on  $[0, 1]$ . To compute these, we can use the `legendreP` command in MATLAB<sup>1</sup>. If we define the function handle

```
lpoly = @(p,z) sqrt(2)*sqrt((2*p+1)/2)*legendreP(p, 2*z-1);
```

Then, for example, `lpoly(3,T)` will return samples of the 3rd order Legendre polynomial at locations in `T`.

- (a) Find the best cubic fit to the data using least-squares. That is, find  $w_0, \dots, w_3$  that minimizes

$$\underset{\mathbf{w}}{\text{minimize}} \quad \sum_{m=1}^M (y_m - f(t_m))^2 \quad \text{where} \quad f(t) = \sum_{n=0}^3 w_n v_n(t),$$

where  $v_n(t)$  is the  $n$ th order Legendre polynomial adapted to  $[0, 1]$ . Let  $\hat{\mathbf{w}}$  be the solution to the above, and  $\hat{f}$  the corresponding cubic polynomial. Compute the sample error

$$\text{sample error} = \left( \sum_{m=1}^M (y_m - \hat{f}(t_m))^2 \right)^{1/2} = \|\mathbf{y} - \mathbf{A}\hat{\mathbf{w}}\|_2,$$

where  $\mathbf{A}$  is the matrix you set up to solve the least-squares problem. Plot your solution  $\hat{f}(t)$  for  $t \in [0, 1]$ , and overlay the sample values  $(t_m, y_m)$  — **the sample values should not have lines connecting them**<sup>2</sup>.

---

<sup>1</sup>For Python, see <https://docs.scipy.org/doc/scipy/reference/generated/scipy.special.legendre.html>.

<sup>2</sup>Use `plot(t,y,'o')` in MATLAB, for example.

- (b) Compute the generalization error

$$\text{generalization error} = \left( \int_0^1 |\hat{f}(t) - f_{\text{true}}(t)|^2 dt \right)^{1/2}.$$

You can either use some numerical integration package to do this (`integral()` in MATLAB), or you can simply sample the functions at 5000 points<sup>3</sup>, take the sum of the squared difference, divide by 5000, then take the square root.

- (c) Repeat part (a) and (b) for polynomials of order  $p = 5, 10, 15, 20, 25$  (note that the number of basis functions is always  $N = p + 1$ ). For each experiment, report the sample error, generalization error, and the largest and smallest singular value of the matrix  $\mathbf{A}$ . Make a plot of the sample error versus  $N$  and the generalization error versus  $N$ . For what value of  $N$  does least-squares start to fall apart and why? Why does the sample error go down monotonically with  $N$  but the generalization error does not? Answer these questions, and for each value of  $p$ , make a plot of the sample values,  $\hat{f}(t)$ , and  $f_{\text{true}}(t)$  overlaid on one another.
- (d) Set  $N = 25$ . Plot the singular values of  $\mathbf{A}$  for this  $N$ . Stabilize your least-squares solution by truncating the SVD; make a judgement call about what  $R'$  should be given your singular value plot, and explain your reasoning. Compute the sample and generalization errors.
- (e) Again with  $N = 25$ , repeat part (d) and compute the truncated SVD estimate for  $R' = 5, 6, \dots, 25$ . Plot the sample and generalization error versus  $R'$ , and interpret in terms of noise error, approximation error, and null space error.
- (f) Fix  $N = 25$ . Now we will consider the ridge regression estimate

$$\underset{\mathbf{w}}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{w}\|_2^2 + \delta \|\mathbf{w}\|_2^2$$

Again examining the plot of the singular values of  $\mathbf{A}$ , come up with a reasonable value of  $\delta$  to use above. Perform the regression, make a plot of your result, again overlaid on the samples and  $f_{\text{true}}(t)$ , and report the sample and generalization error. Also comment on what happens to both the sample and generalization error as you sweep the value of  $\delta$  — provide representative plots.

2. Let

$$\mathbf{H} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ -3 \end{bmatrix},$$

and let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

- (a) What is the smallest value that  $f$  takes on  $\mathbb{R}^2$ ? At what  $\mathbf{x}$  does it achieve this minimum value?
- (b) Write  $f(\mathbf{x})$  out as a quadratic function in  $x_1, x_2$  where  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ . In other words, fill in the blanks below

$$f(\mathbf{x}) = \underline{\hspace{1cm}} x_1^2 + \underline{\hspace{1cm}} x_2^2 + \underline{\hspace{1cm}} x_1 x_2 + \underline{\hspace{1cm}} x_1 + \underline{\hspace{1cm}} x_2.$$

<sup>3</sup>5000 might be overkill here, but these computations are cheap.

- (c) Using MATLAB or Python, make a contour plot of  $f(\mathbf{x})$  around its minimizer in  $\mathbb{R}^2$ . Compute the eigenvectors and eigenvalues of  $\mathbf{H}$ , and discuss what role they are playing in the geometry of your sketch.
  - (d) On top of the contour plot, trace out the first four steps of the gradient descent algorithm starting at  $\mathbf{x}_0 = \mathbf{0}$ .
  - (e) On top of the contour plot, trace out the two steps of the conjugate gradient method starting at  $\mathbf{x} = \mathbf{0}$ .
3. (a) Write a MATLAB (or Python) function `gdsolve` that implements the gradient descent algorithm for solving linear systems of equations. The function should be called as

```
[x, iter] = gdsolve(H, b, tol, maxiter);
```

where  $\mathbf{H}$  is a  $N \times N$  symmetric positive definite matrix,  $\mathbf{b}$  is a vector of length  $N$ , and `tol` and `maxiter` specify the halting conditions. Your algorithm should iterate until  $\|\mathbf{b} - \mathbf{H}\mathbf{x}_k\|_2 / \|\mathbf{b}\|_2$  is less than `tol` or the maximum number of iterations `maxiter` has been reached. For the outputs: `x` is your solution, and `iter` is the number of iterations that were performed. Run your code on the  $\mathbf{H}$  and  $\mathbf{b}$  in the file `hw05p3_data.mat` for a `tol` of  $10^{-6}$ . Report the number of iterations needed for convergence, and for your solution  $\hat{\mathbf{x}}$  verify that  $\|\mathbf{b} - \mathbf{H}\hat{\mathbf{x}}\|_2 / \|\mathbf{b}\|_2$  is within the specified tolerance.

- (b) Write a MATLAB (or Python) function `cgsolve.m` that implements the conjugate gradient method. The function should be called as

```
[x, iter] = cgsolve(H, b, tol, maxiter);
```

where the inputs and outputs have the same interpretation as in the previous problem. Run your code on the  $\mathbf{H}$  and  $\mathbf{b}$  in the file `hw05p3_data.mat` (same data as in part(a)) for a `tol` of  $10^{-6}$ . Report the number of iterations needed for convergence, and for your solution  $\hat{\mathbf{x}}$  verify that  $\|\mathbf{b} - \mathbf{H}\hat{\mathbf{x}}\|_2 / \|\mathbf{b}\|_2$  is within the specified tolerance. Compare against the gradient descent algorithm in part (a).

4. Let  $\mathbf{A}$  be a banded tridiagonal matrix:

$$\mathbf{A} = \begin{bmatrix} d_1 & c_1 & 0 & 0 & 0 & \cdots & 0 \\ f_1 & d_2 & c_2 & 0 & 0 & \cdots & 0 \\ 0 & f_2 & d_3 & c_3 & 0 & \cdots & 0 \\ 0 & 0 & f_3 & d_4 & c_4 & \cdots & 0 \\ \vdots & \vdots & & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & & f_{N-2} & d_{N-1} & c_{N-1} \\ 0 & 0 & 0 & \cdots & & f_{N-1} & d_N \end{bmatrix}$$

- (a) Argue that the LU factorization of  $\mathbf{A}$  has the form

$$\mathbf{A} = \begin{bmatrix} * & 0 & 0 & 0 & \cdots & 0 \\ * & * & 0 & 0 & \cdots & 0 \\ 0 & * & * & 0 & \cdots & 0 \\ 0 & 0 & * & * & \cdots & 0 \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & & * & * \end{bmatrix} \begin{bmatrix} * & * & 0 & 0 & 0 & \cdots & 0 \\ 0 & * & * & 0 & 0 & \cdots & 0 \\ 0 & 0 & * & * & 0 & \cdots & 0 \\ \vdots & & & \ddots & \ddots & & \\ 0 & 0 & \cdots & & & * & * \\ 0 & 0 & \cdots & & & 0 & * \end{bmatrix},$$

where  $*$  signifies a non-zero term.

- (b) Write down an algorithm that computes the LU factorization of  $\mathbf{A}$ , meaning the  $\{\ell_i\}$ ,  $\{u_i\}$ , and  $\{g_i\}$  below

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ \ell_1 & 1 & 0 & 0 & \cdots & 0 \\ 0 & \ell_2 & 1 & 0 & \cdots & 0 \\ 0 & 0 & \ell_3 & 1 & \cdots & 0 \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & & \ell_{N-1} & 1 \end{bmatrix} \begin{bmatrix} u_1 & g_1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & u_2 & g_2 & 0 & 0 & \cdots & 0 \\ 0 & 0 & u_3 & g_3 & 0 & \cdots & 0 \\ \vdots & & & \ddots & \ddots & & \\ 0 & 0 & \cdots & & u_{N-1} & g_{N-1} \\ 0 & 0 & \cdots & & 0 & u_N \end{bmatrix},$$

I will get you started:

```

u1 = d1
for k = 2, ..., N
    gk-1 =
    ℓk-1 =
    uk =
end

```

- (c) What is the computational complexity of the algorithm above? (How does the number of computations scale with  $N$ ?) Once the LU factorization is in hand, how does solving  $\mathbf{Ax} = \mathbf{b}$  scale with  $N$ ?
5. The Gram-Schmidt process is a method for orthonormalizing a set of vectors in an inner product space. The method works as follows: if  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N\}$  is a set of linearly independent vectors in  $\mathbb{R}^M$  (so clearly  $N \leq M$ ) then we can generate a sequence of orthonormal vectors  $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N\}$  such that

$$\text{Span}(\{\mathbf{a}_1, \dots, \mathbf{a}_N\}) = \text{Span}(\{\mathbf{q}_1, \dots, \mathbf{q}_N\})$$

using

$$\mathbf{q}_1 = \frac{\mathbf{a}_1}{\|\mathbf{a}_1\|_2}$$

and for  $k = 2, \dots, N$ :

$$\mathbf{w}_k = \mathbf{a}_k - \sum_{\ell=1}^{k-1} \langle \mathbf{a}_k, \mathbf{q}_\ell \rangle \mathbf{q}_\ell,$$

$$\mathbf{q}_k = \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|_2}.$$

- (a) As a warm-up, find  $\mathbf{q}_1, \mathbf{q}_2$  and  $\mathbf{q}_3$  when

$$\mathbf{a}_1 = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix}, \quad \mathbf{a}_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{a}_3 = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \\ 1 \end{bmatrix}.$$

Feel free to use a computer to do the calculations; just explain what you did in your write-up.

- (b) For  $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3\}$  and  $\{\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3\}$  as in part (a), let

$$\mathbf{A} = \begin{bmatrix} | & | & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 \\ | & | & | \end{bmatrix}, \quad \mathbf{Q} = \begin{bmatrix} | & | & | \\ \mathbf{q}_1 & \mathbf{q}_2 & \mathbf{q}_3 \\ | & | & | \end{bmatrix}.$$

Show how we can write  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{R}$  is upper triangular. Do this by explicitly calculating  $\mathbf{R}$ . (Hint: just keep track of your work while doing part (a)).

- (c) Suppose I run the algorithm above on a general  $M \times N$  matrix  $\mathbf{A}$  with linearly independent columns (full column rank). Explain how the Gram-Schmidt algorithm can be interpreted as finding a  $M \times N$  matrix  $\mathbf{Q}$  with orthonormal columns and an upper triangular matrix  $\mathbf{R}$  such that  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ . Do this by explicitly writing what the entries of  $\mathbf{R}$  are in terms of the quantities that appear in the algorithm. This is called the *QR decomposition* of  $\mathbf{A}$ .
- (d) Prove that an upper triangular matrix is invertible if and only if the elements along the diagonal are nonzero. Using this to show that the linear independence of the columns of  $\mathbf{A}$  implies that all of the entries along the diagonal of  $\mathbf{R}$  will be nonzero?
- (e) Suppose that an  $M \times N$  matrix  $\mathbf{A}$  has full column rank. Show that the solution to the least-squares problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

is  $\hat{\mathbf{x}} = \mathbf{R}^{-1}\mathbf{Q}^T\mathbf{y}$ , where  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  is the QR decomposition of  $\mathbf{A}$ .