

Math Foundations of ML, Fall 2018

Homework #6

Due Friday October 19, at the beginning of class

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

1. Using your class notes, prepare a 1-2 paragraph summary of what we talked about in class in the last week. I do not want just a bulleted list of topics, I want you to use complete sentences and establish context (Why is what we have learned relevant? How does it connect with other things you have learned here or in other classes?). The more insight you give, the better.
2. Recall the Gram-Schmidt algorithm from HW 3: 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

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

then for $k = 2, 3, \dots$,

$$\mathbf{w}_k = \mathbf{a}_k - \sum_{\ell=1}^{k-1} \langle \mathbf{a}_k, \mathbf{q}_\ell \rangle \mathbf{q}_\ell$$
$$\mathbf{q}_k = \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|_2}.$$

- (a) As a warm-up and review, find $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3$ 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 the last part, 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 $\mathbf{A} = \mathbf{QR}$, 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 $\mathbf{A} = \mathbf{QR}$. 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 or disprove: an upper triangular matrix is invertible if and only none of the elements along the diagonal are zero. Why does the linear independence of the columns of \mathbf{A} mean 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 $\hat{x} = R^{-1}Q^T y$, where $A = QR$ is the QR decomposition of A .

3. As we have seen (i.e. Homework 1, Problem 5), a function $f : \mathbb{R} \rightarrow \mathbb{R}$ of one variable that is twice differentiable is convex if its second derivative is positive everywhere. A function $f : \mathbb{R}^M \rightarrow \mathbb{R}$ of M variables is convex if every function of one variable formed by looking along a ray starting from an arbitrary point has a second derivative that is positive anywhere. That is, the function

$$g_{\mathbf{x}, \mathbf{v}}(t) = f(\mathbf{x} + t\mathbf{v})$$

is a convex function of one variable (t) for all \mathbf{x}, \mathbf{v} .

Let \mathbf{H} be a symmetric matrix. Show that

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{H} \mathbf{x}$$

is convex if and only if \mathbf{H} is symmetric positive semi-definite. That is, all the eigenvalues of \mathbf{H} are non-negative, or equivalently $\mathbf{w}^\top \mathbf{H} \mathbf{w} \geq 0$ for all $\mathbf{w} \in \mathbb{R}^M$.

4. 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

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

- What is the smallest value that f takes on \mathbb{R}^2 ? At what \mathbf{x}_\star does it achieve this minimum value?
- 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

$$f(\boldsymbol{x}) = __x_1^2 + __x_2^2 + __x_1x_2 + __x_1 + __x_2.$$

- (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 $\mathbf{x}_0 = \mathbf{0}$.
5. A look at the plots on pages 119 and 123 as well as the plot you made in answering the previous question should give you the following intuition: if the initial point \mathbf{x}_0 is along one of the axes of the ellipses defined by the contour lines of $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{x}^T \mathbf{b}$, then steepest descent converges in one step, i.e. $\mathbf{x}_1 = \mathbf{x}_*$.

We will now make this intuition formal. Suppose

$$\mathbf{x}_0 = \mathbf{x}_* + \beta \mathbf{v},$$

where \mathbf{v} is an eigenvector of \mathbf{H} with associated eigenvalue λ . Show that

$$\mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{r}_0 = \mathbf{x}_*,$$

where $\mathbf{r}_0 = \mathbf{b} - \mathbf{H}\mathbf{x}_0$ and α_0 is chosen as on page 121 of the notes.

6. Write a MATLAB function `sdsolve` that implements the steepest descent algorithm. The function should be called as

```
[x, iter] = sdsolve(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 $\|\mathbf{r}_k\|_2 / \|\mathbf{b}\|_2$ 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 `hw6p6.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 $\mathbf{H}\hat{\mathbf{x}}$ is within the specified tolerance of \mathbf{b} .