

# MATH661 Homework 3 - Least squares problems

BY LEYI ZHANG

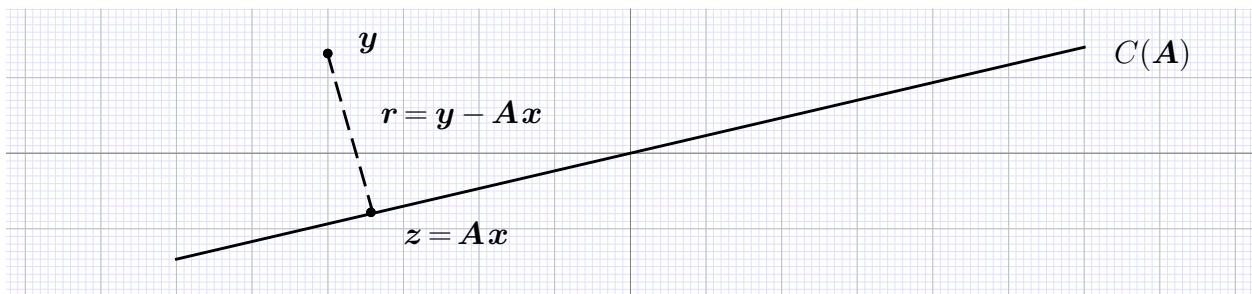
09/15/21

This assignment addresses one of the fundamental topics within scientific computation: finding economical descriptions of complex objects. Some object is described by  $\mathbf{y} \in \mathbb{C}^m$  (with  $m$  typically large), and a reduced description is sought by linear combination  $\mathbf{A}\mathbf{x}$ , with  $\mathbf{A} \in \mathbb{C}^{m \times n}$  ( $n < m$ , often  $n \ll m$ ). The surprisingly simple Euclidean geometry of Fig. 1 (which should be committed to memory) will be shown to have wide-ranging applicability to many different types of problems. The error (or residual) in approximating  $\mathbf{y}$  by  $\mathbf{A}\mathbf{x}$  is defined as

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x},$$

and 2-norm minimization defines the least-squares problem

$$\min_{\mathbf{x} \in \mathbb{C}^m} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|.$$



**Figure 1.** Least squares (2-norm error minimization) problem.

## Track 1

Consider data  $\mathcal{D} = \{(t_i, y_i) | i = 1, 2, \dots, m\}$  obtained by sampling a function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , with  $y_i = f(t_i)$ . An approximation is sought by linear combination

$$f(t) \cong x_1 a_1(t) + x_2 a_2(t) + \dots + x_n a_n(t).$$

Introduce the vector-valued function  $A: \mathbb{R} \rightarrow \mathbb{R}^n$  (organized as a row vector)

$$A(t) = [ a_1(t) \ a_2(t) \ \dots \ a_n(t) ],$$

such that

$$f(t) \cong A(t) \mathbf{x}, \mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T.$$

With  $\mathbf{t} = [t_1 \ t_2 \ \dots \ t_m]^T$  a sampling of the function domain, a matrix is defined by

$$\mathbf{A} = A(\mathbf{t}) \mathbf{x} = [a_1(\mathbf{t}) \ a_2(\mathbf{t}) \ \dots \ a_n(\mathbf{t})] \mathbf{x} \in \mathbb{R}^{m \times n}.$$

**Tasks.** In each exercise below, construct the least-squares approximant for the stated range of  $n \in \mathcal{N}$ , sample points  $\mathbf{t}$ , and choice of  $A(t)$ . Plot in a single figure all components of  $A(t)$ . Plot the approximants, as well as  $f$  in a single figure. Construct a convergence plot of the approximations by representation of point data  $\mathcal{E} = \{(\log n, \log \|\mathbf{y} - \mathbf{A}\mathbf{x}\|) \mid \mathbf{A} \in \mathbb{R}^{m \times n}, n \in \mathcal{N}\}$ . For the largest value of  $n$  within  $\mathcal{N}$ , construct a figure superimposing increasing number of sampling points,  $m \in \mathcal{M}$ . Comment on what you observe in each individual exercise. Also compare results from the different exercises.

### Problem 1.1

Start with the classical example due to Runge (1901)

$$\begin{aligned} f: [-1, 1] &\rightarrow \mathbb{R}, f(t) = \frac{1}{(1 + 25t^2)}, t_i = \frac{2(i-1)}{m-1} - 1, \\ \mathcal{M} &= \{16, 32, 64, 128, 256\}, \mathcal{N} = \{4, 8, 16, 32\}, \\ A(t) &= [1 \ t \ t^2 \ \dots \ t^{n-1}]. \end{aligned}$$

### Problem 1.2

Instead of the equidistant point samples of the Runge example above use the Chebyshev nodes

$$t_i = \cos\left(\frac{2i-1}{2m}\pi\right),$$

keeping other parameters as in Problem 1.

### Problem 1.3

Instead of the monomial family of the Runge example, use the Fourier basis

$$A(t) = [1 \ \cos \pi t \ \sin \pi t \ \dots \ \cos \pi n t \ \sin \pi n t]$$

keeping other parameters as in Problem 1. In this case  $\mathbf{A} \in \mathbb{R}^{m \times (2n+1)}$ .

## Problem 1.4

Instead of the monomial family of the Runge example, use the piecewise linear  $B$ -spline basis

$$A(t) = [N_1(t) \ N_2(t) \ \dots \ N_n(t)],$$

$$N_i(t) = \begin{cases} 0, & t < t_{i-1} \\ \frac{t - t_{i-1}}{h} & t_{i-1} \leq t < t_i \\ \frac{t_{i+1} - t}{h} & t_i \leq t < t_{i+1} \\ 0 & t_{i+1} \leq t \end{cases}, h = \frac{2}{m-1},$$

keeping other parameters as in Problem 1.

```
∴ using Plots; using LinearAlgebra
```

```
∴
```

## Track 2

### Problem 2.1

If  $\mathbf{Q} \in \mathbb{C}^{m \times n}$  has orthonormal columns, prove that  $\mathbf{P}_\mathbf{Q} = \mathbf{Q}\mathbf{Q}^*$  is an orthogonal projector onto  $C(\mathbf{Q})$ . Determine the expression of  $\mathbf{P}_\mathbf{A}$ , the projector onto  $C(\mathbf{A})$ , with  $\mathbf{A} \in \mathbb{C}^{m \times n}$ . Compare the number of arithmetic operations required to compute  $\mathbf{y} = \mathbf{P}_\mathbf{A} \mathbf{x}$ , by comparison to first determining the  $QR$  factorization,  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , and then computing  $\mathbf{y} = \mathbf{Q}\mathbf{Q}^* \mathbf{x}$ .

**Solution:**

**Definition 1.** A orthogonal projector is a square matrix  $\mathbf{P}$  that satisfies  $\mathbf{P}^2 = \mathbf{P}$  and  $\mathbf{P}^* = \mathbf{P}$ .

Let  $\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_n] \in \mathbb{C}^{m \times n}$ , where  $\mathbf{q}_j = [q_{j,1} \ q_{j,2} \ \dots \ q_{j,m}]^T$  is a column vector with  $q_{j,i} \in \mathbb{C}$  for every  $i \in \{1, \dots, m\}$  and  $j \in \{1, \dots, n\}$ . Then  $\mathbf{Q}^* = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_n]^* \in \mathbb{C}^{n \times m}$ . If we write them out, we have

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_n] = \begin{bmatrix} q_{1,1} & q_{2,1} & \dots & q_{n,1} \\ q_{1,2} & q_{2,2} & & \vdots \\ \vdots & & \ddots & \vdots \\ q_{1,m} & \dots & \dots & q_{n,m} \end{bmatrix},$$

$$\mathbf{Q}^* = \begin{bmatrix} \mathbf{q}_1^* \\ \mathbf{q}_2^* \\ \vdots \\ \mathbf{q}_n^* \end{bmatrix} = \begin{bmatrix} \overline{q_{1,1}} & \overline{q_{1,2}} & \dots & \overline{q_{1,m}} \\ \overline{q_{2,1}} & \overline{q_{2,2}} & & \vdots \\ \vdots & & \ddots & \vdots \\ \overline{q_{n,1}} & \dots & \dots & \overline{q_{n,m}} \end{bmatrix}.$$

First, we want to show that  $P_Q^* = P_Q$ . We know that for any matrices  $A, B$ , we have  $(AB)^* = B^*A^*$ . Thus

$$P_Q^* = (QQ^*)^* = (Q^*)^*Q^* = QQ^* = P_Q.$$

Then we want to show that  $P_Q^2 = P_Q$ . We are given that  $P_Q = QQ^*$ , so we have

$$P_Q^2 = QQ^*QQ^*.$$

Consider

$$Q^*Q = \begin{bmatrix} q_1^* \\ q_2^* \\ \vdots \\ q_n^* \end{bmatrix} [q_1 \ q_2 \ \cdots \ q_n] = \begin{bmatrix} q_1^*q_1 & q_1^*q_2 & \cdots & q_1^*q_n \\ q_2^*q_1 & q_2^*q_2 & & \vdots \\ \vdots & & \ddots & \vdots \\ q_n^*q_1 & \cdots & \cdots & q_n^*q_n \end{bmatrix} \in \mathbb{C}^{n \times n}.$$

Since  $Q$  has orthonormal columns, we know that  $q_i^*q_j = \delta_{ij}$  for every  $i, j \in \{1, \dots, n\}$ . That gives us

$$Q^*Q = I_n,$$

and we have

$$P_Q^2 = QQ^*QQ^* = Q(Q^*Q)Q^* = QI_nQ^* = QQ^* = P_Q.$$

Therefore, we conclude that  $P_Q$  is indeed an orthogonal projector.

Suppose we have a matrix  $A = [a_1 \ a_2 \ \cdots \ a_n] \in \mathbb{C}^{m \times n}$ , where  $a_j = [a_{j,1} \ a_{j,2} \ \cdots \ a_{j,m}]^T$  with  $a_{j,i} \in \mathbb{C}$  for every  $i \in \{1, \dots, m\}$  and  $j \in \{1, \dots, n\}$ . Then  $A^* = [a_1 \ a_2 \ \cdots \ a_n]^* \in \mathbb{C}^{n \times m}$ . That means  $P_A = AA^* \in \mathbb{C}^{m \times m}$  such that

$$\begin{aligned} P_A = AA^* &= \begin{bmatrix} \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,1} & \sum_{j=1}^n \overline{a_{j,2}} \cdot a_{j,1} & \cdots & \sum_{j=1}^n \overline{a_{j,m}} \cdot a_{j,1} \\ \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,2} & \sum_{j=1}^n \overline{a_{j,2}} \cdot a_{j,2} & & \vdots \\ \vdots & & \ddots & \vdots \\ \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,m} & \cdots & \cdots & \sum_{j=1}^n \overline{a_{j,m}} \cdot a_{j,m} \end{bmatrix}. \\ P_A = AA^* &= \begin{bmatrix} \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,1} & \sum_{j=1}^n \overline{a_{j,2}} \cdot a_{j,1} & \cdots & \sum_{j=1}^n \overline{a_{j,m}} \cdot a_{j,1} \\ \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,2} & \sum_{j=1}^n \overline{a_{j,2}} \cdot a_{j,2} & & \vdots \\ \vdots & & \ddots & \vdots \\ \sum_{j=1}^n \overline{a_{j,1}} \cdot a_{j,m} & \cdots & \cdots & \sum_{j=1}^n \overline{a_{j,m}} \cdot a_{j,m} \end{bmatrix} \\ &= \begin{bmatrix} \sum_{j=1}^n \overline{a_{j,1}} a_j & \sum_{j=1}^n \overline{a_{j,2}} a_j & \cdots & \sum_{j=1}^n \overline{a_{j,m}} a_j \end{bmatrix}. \end{aligned}$$

For every entry of  $\mathbf{P}_A$ , we will need to carry out  $n$  multiplications and  $n - 1$  additions, i.e.  $2n - 1$  operations in total. Since  $\mathbf{P}_A$  is a  $m \times m$  matrix, it has  $m^2$  entries. Thus we will need to carry out  $m^2(2n - 1)$  arithmetic operations to compute  $\mathbf{P}_A$ . Now suppose  $\mathbf{x} \in \mathbb{C}^m$ , then we need another  $m$  multiplications and  $m - 1$  additions to compute each of the  $m$  entries of  $\mathbf{P}_A \mathbf{x}$ . Therefore, the total number of operations to compute  $\mathbf{P}_A \mathbf{x}$  will be

$$m^2(2n - 1) + m(2m - 1) = m^2(2n + 1) - m.$$

Assume we are using the modified Gram-Schmidt process to carry out the QR decomposition of  $A$ . Consider  $\mathbf{q}_i = \mathbf{a}_i$  with length  $m$ , one of the  $n$  columns of  $A$ . We first need to compute its norm  $\mathbf{q}_i^* \mathbf{q}_i$ , which will take  $m$  multiplications,  $m - 1$  additions, and  $m$  square roots, so  $3m - 1$  operations. Then we normalize  $\mathbf{q}_i$  by dividing each entry by its norm, and that takes  $m$  operations. Then for each  $j \in \{i + 1, \dots, n\}$ , we will compute  $\mathbf{R}_{ij}$  with a dot product  $\mathbf{q}_i^* \mathbf{a}_j$  with  $m$  multiplications and  $m - 1$  additions, and then compute  $\mathbf{q}_j$  with a dot product  $\mathbf{r}_i^* \mathbf{q}_j$  with  $m$  multiplications and  $m - 1$  additions, and then subtract the dot product from each entry of  $\mathbf{q}_j$  with  $m$  operations. That is, for each  $j$ , we need  $(2m - 1) + (3m - 1) = 5m - 2$  operations, and there are  $n - i - 1$  of  $j$ 's for each  $i$  from 1 to  $n - 1$ . Thus for each  $i$  from 1 to  $n - 1$ , we need  $(4m - 1) + (5m - 2)(n - i - 1)$  operations, and for  $i = n$ , only  $4m - 1$  operations are needed.

Therefore, the total number of operations we will need for the modified Gram-Schmidt process is

$$(4m - 1)n + \sum_{i=1}^{n-1} (5m - 2)(n - i - 1) = (4m - 1)n + (5m - 2)(n - 2)(n - 1).$$

Now we know that  $\mathbf{Q}^* \in \mathbb{C}^{n \times m}$ ,  $\mathbf{Q} \in \mathbb{C}^{m \times n}$ , and  $\mathbf{x} \in \mathbb{C}^m$ . So each of the  $n$  entries of  $\mathbf{Q}^* \mathbf{x}$  will take  $m$  multiplications and  $m - 1$  additions to compute, so  $n(2m - 1)$  in total for  $\mathbf{Q}^* \mathbf{x}$ . Then for each of the  $m$  entries of  $\mathbf{Q} \mathbf{Q}^* \mathbf{x}$ , we will need  $n$  multiplications and  $n - 1$  additions, so  $m(2n - 1)$  in total for  $\mathbf{Q} \mathbf{Q}^* \mathbf{x}$ . Thus we need  $n(2m - 1) + m(2n - 1) = 4mn - m - n$  operations to compute  $\mathbf{Q} \mathbf{Q}^* \mathbf{x}$ .

In the end, the number of operations we need to first determine the QR decomposition of  $A$  and then compute  $\mathbf{Q} \mathbf{Q}^* \mathbf{x}$  would be

$$(5m - 2)n^2 - 14mn + 9m + 4n - 4.$$

Thus the cost of computing  $\mathbf{P}_A \mathbf{x}$  directly would be of order  $m^2n$ , while the cost of computing  $\mathbf{Q} \mathbf{Q}^* \mathbf{x}$  after QR decomposition would be of order  $mn^2$ . The QR decomposition method will be much more efficient in the situations where we have  $m \gg n$ , which is often the case.

## Problem 2.2

Continuing Problem 1, determine  $\|\mathbf{P}_Q\|_2$ , and express  $\|\mathbf{P}_A\|_2$  in terms of the singular value decomposition of  $A$ . Comment the result, considering, say, length of shadows at various times of day.

**Solution:**

**Definition 2.** The 2-norm of a matrix  $A \in \mathbb{C}^{m \times n}$  is defined to be

$$\|A\|_2 = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2=1} \|A\mathbf{x}\|_2.$$

Since  $P_Q = P_Q^*$  and  $P_Q = P_Q^2$ , for any  $\mathbf{x} \in \mathbb{C}^m$ , we have

$$\|P_Q \mathbf{x}\|_2^2 = (P_Q \mathbf{x})^* P_Q \mathbf{x} = \mathbf{x}^* P_Q^* P_Q \mathbf{x} = \mathbf{x}^* P_Q^2 \mathbf{x} = \mathbf{x}^* P_Q \mathbf{x}.$$

Consider  $Q: \mathbb{C}^n \rightarrow \mathbb{C}^m$ . We know that the codomain  $\mathbb{C}^m = C(Q) \oplus N(Q^*)$ . That means any vector  $\mathbf{x} \in \mathbb{C}^m$ , can be written as  $\mathbf{x} = \mathbf{v} + \mathbf{w}$  with unique vectors  $\mathbf{v} \in C(Q)$  and  $\mathbf{w} \in N(Q^*)$ .

Since  $\mathbf{v}$  is an element in  $C(Q)$ , there exists some  $\mathbf{y} \in \mathbb{C}^n$ , such that  $\mathbf{v} = Q\mathbf{y}$ . Since  $P_Q = QQ^*$ , we get  $P_Q \mathbf{v} = QQ^*Q\mathbf{y}$ . We have shown earlier that  $Q^*Q = I_n$ , so  $P_Q \mathbf{v} = Q\mathbf{y} = \mathbf{v}$ .

Since  $\mathbf{w}$  is an element in  $N(Q^*)$ , we have  $Q^*\mathbf{w} = 0$  by definition. Thus we get  $P_Q \mathbf{w} = QQ^*\mathbf{w} = 0$ .

Therefore, we have

$$\begin{aligned} \|P_Q \mathbf{x}\|_2^2 &= \mathbf{x}^* P_Q \mathbf{x} = (\mathbf{v} + \mathbf{w})^* P_Q (\mathbf{v} + \mathbf{w}) \\ &= (\mathbf{v}^* + \mathbf{w}^*) (P_Q \mathbf{v} + P_Q \mathbf{w}) \\ &= (\mathbf{v}^* + \mathbf{w}^*) \mathbf{v} \\ &= \mathbf{v}^* \mathbf{v} + \mathbf{w}^* \mathbf{v} \end{aligned}$$

We know that  $C(Q) \perp N(Q^*)$ , so we have  $\mathbf{w}^* \mathbf{v} = 0$ , and

$$\|P_Q \mathbf{x}\|_2^2 = \mathbf{v}^* \mathbf{v} = \|\mathbf{v}\|_2^2.$$

Since  $\mathbf{x} = \mathbf{v} + \mathbf{w}$  and  $\|\mathbf{x}\|_2 = 1$ ,

$$\begin{aligned} 1 = \|\mathbf{x}\|_2^2 &= (\mathbf{v} + \mathbf{w})^* (\mathbf{v} + \mathbf{w}) \\ &= (\mathbf{v}^* + \mathbf{w}^*) (\mathbf{v} + \mathbf{w}) \\ &= \mathbf{v}^* \mathbf{v} + \mathbf{v}^* \mathbf{w} + \mathbf{w}^* \mathbf{v} + \mathbf{w}^* \mathbf{w} \\ &= \|\mathbf{v}\|_2^2 + \|\mathbf{w}\|_2^2. \end{aligned}$$

That means we have  $\|\mathbf{v}\|_2^2 \leq 1$  with  $\|\mathbf{v}\|_2^2 = 1$  when  $\|\mathbf{w}\|_2^2 = 0$ . Thus we conclude that

$$\sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2=1} \|P_Q \mathbf{x}\|_2^2 = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2=1} \|\mathbf{v}\|_2^2 = 1,$$

and

$$\|P_Q\|_2 = \sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2=1} \|P_Q \mathbf{x}\|_2 = \sqrt{\sup_{\mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2=1} \|P_Q \mathbf{x}\|_2^2} = 1.$$

Suppose there exist matrices  $\mathbf{U} \in \mathbb{C}^{m \times m}$ ,  $\mathbf{V} \in \mathbb{C}^{n \times n}$  and  $\mathbf{\Sigma} \in \mathbb{C}^{m \times n}$  with  $\mathbf{U}^* \mathbf{U} = \mathbf{I}_m$  and  $\mathbf{V}^* \mathbf{V} = \mathbf{I}_n$ , such that

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*.$$

Then

$$\begin{aligned} \mathbf{A} \mathbf{A}^* &= (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^*) (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^*)^* \\ &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* \mathbf{V} \mathbf{\Sigma}^* \mathbf{U}^* \\ &= \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^* \mathbf{U}^*, \end{aligned}$$

and we have

$$\begin{aligned} \|\mathbf{P}_\mathbf{A} \mathbf{x}\|_2^2 &= (\mathbf{P}_\mathbf{A} \mathbf{x})^* \mathbf{P}_\mathbf{A} \mathbf{x} = \mathbf{x}^* \mathbf{P}_\mathbf{A}^* \mathbf{P}_\mathbf{A} \mathbf{x} = \mathbf{x}^* \mathbf{P}_\mathbf{A}^2 \mathbf{x} \\ &= \mathbf{x}^* \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^* \mathbf{U}^* \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^* \mathbf{U}^* \mathbf{x} \\ &= \mathbf{x}^* \mathbf{U} \mathbf{\Sigma} \mathbf{\Sigma}^* \mathbf{\Sigma} \mathbf{\Sigma}^* \mathbf{U}^* \mathbf{x} \\ &= \mathbf{x}^* \mathbf{U} (\mathbf{\Sigma} \mathbf{\Sigma}^*)^2 \mathbf{U}^* \mathbf{x}. \end{aligned}$$

In this problem, we found that the projection, or the “shadow”, of a vector on the the comlumn space of a matrix cannot be longer than itself, but in real life the length of the shadow can have arbitrary length depending on the angle of the incoming light. I guess in this case, the projection is more similar to the shadow created by a light source that is directly above the object, like the shadow of an object on the equator at noon.

### Problem 2.3

A matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{C}^{m \times n}$  is said to be banded with bandwidth  $B$  if  $a_{ij} = 0$  for  $|i - j| > B$ . Implement the modified Gram-Schmidt algorithm for  $\mathbf{A} \in \mathbb{C}^{m \times n}$  a banded matrix with bandwidth  $B$  using as few arithmetic operations as possible.

(I don't have a working solution for this problem yet. I will read more and work on it over the next week.)

If  $B < 0$ , then  $\mathbf{A}$  would be a matrix of zeros, so we only consider the case where  $B \geq 0$ . Let  $b \in \mathbb{N}$  be the largest natural number such that  $b \leq B$ , i.e. we have  $b \leq B < b + 1$ . If  $|i - j| > B$  for any  $i, j$ , then  $|i - j| \geq b + 1$ .

Suppose  $\mathbf{A} = [\mathbf{a}_1 \mathbf{a}_2 \cdots \mathbf{a}_n]$  where  $\mathbf{a}_j = [a_{j,1} \ a_{j,2} \cdots a_{j,m}]$  with  $a_{j,k} \in \mathbb{C}$  for all  $j \in \{1, \dots, n\}$  and  $k \in \{1, \dots, m\}$ . Then for each  $j$ , we have  $a_{j,1} = \dots = a_{j,j-b-1} = a_{j,j+b+1} = \dots = a_{j,m} = 0$ , i.e. there are at most  $2b + 1$  nonzero entries in each  $\mathbf{a}_j$ .

For any  $\mathbf{a}_j$  with  $j < b + 1$ , we only need to consider  $j + b$  entries. For any  $\mathbf{a}_j$  with  $j > m - b$ , we need to consider  $m + b - j$  entries. For  $\mathbf{a}_j$  with  $b + 1 \leq j \leq m - b$ , we need to consider  $2b + 1$  entries.

## Problem 2.4

(Solve Problem 1, Track 1.) Start with the classical example due to Runge (1901)

$$f: [-1, 1] \rightarrow \mathbb{R}, f(t) = \frac{1}{(1 + 25t^2)}, t_i = \frac{2(i-1)}{m-1} - 1, \\ \mathcal{M} = \{16, 32, 64, 128, 256\}, \mathcal{N} = \{4, 8, 16, 32\}, \\ A(t) = \begin{bmatrix} 1 & t & t^2 & \dots & t^{n-1} \end{bmatrix}.$$

Julia (1.6.2) session in GNU TeXmacs

```
∴ function runge_f(t)
    f_t = 1 ./ (25*(t.^2) .+1)
    return f_t
end
```

runge\_f

```
∴ function monomial(n,m)
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = ones(m,1)
    for i=1:n-1
        A = [A (t.^ i)]
    end
    return A
end
```

monomial

```
∴ function monomial_coef(n,m)
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = monomial(n,m)
    y = runge_f.(t)
    x = A\y
    return x
end
```

monomial\_coef

```
∴ M = 512; H = 2.0/(M-1); T = [(i*H-1) for i=0:(M-1)]; Y = runge_f(T);
```

```
∴ m_values = 2.^ collect(4:8); m_size = length(m_values); err =
    Array{Float64}(undef, m_size);
```

```
∴
```



```

∴ for i=1:4
    n = 2^(i+1)
    m=96
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = ones(m,1)
    for i=1:n-1
        A = [A (t .^ i)]
    end
    display(Plots.plot(A, legend=false))
end

```

⌈ Warning: Assignment to `n` in soft scope is ambiguous because a global variable by the same name exists: `n` will be treated as a new local. Disambiguate by using `local n` to suppress this warning or `global n` to assign to the existing global variable.  
 ⌋ @ none:2

```
Base.Meta.ParseError("extra_token_after_end_of_expression")
```

```
∴
```

```
∴ n=4; B=monomial(n,M); Approx=zeros(M,m_size);
```

```

∴ for i=1:m_size
    m=m_values[i];
    x=monomial_coef(n,m);
    Approx[:,i]=B*x;
end

```

```
∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"));
```

```
∴
```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

```

```
∴ err = log.(err .^0.5);
```

```
∴ display(Plots.plot(log.(m_values),err, label="error"))
```

```

∴ num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

```

```
∴ print("The_order_of_convergence_is_estimated_to_be","\n",q)
```

The order of convergence is estimated to be  
 0.2464832341748203

```
∴ print("The rate of convergence is", "\n", s)
```

The rate of convergence is  
1.1570229053972054

```
∴
```

```
∴ n=8; B=monomial(n,M); Approx=zeros(M,m_size);
```

```
∴ for i=1:m_size  
    m=m_values[i];  
    x=monomial_coef(n,m);  
    Approx[:,i]=B*x;  
end
```

```
∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"));
```

```
∴
```

```
∴ for i=1:m_size  
    err_abs = abs.(Y-Approx[:,i])  
    err[i] = dot(err_abs,err_abs)  
end
```

```
∴ err = log.(err.^0.5);
```

```
∴ display(Plots.plot(log.(m_values),err, label="error"))
```

```
∴ num=err[3:m_size] .- err[2:(m_size-1)];  
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);  
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));
```

```
∴ print("The order of convergence is estimated to be", "\n", q)
```

The order of convergence is estimated to be  
0.40277320620291507

```
∴ print("The rate of convergence is", "\n", s)
```

The rate of convergence is  
1.5031719274633308

```
∴
```

```
∴ n=16; B=monomial(n,M); Approx=zeros(M,m_size);
```

```
∴ for i=1:m_size  
    m=m_values[i];  
    x=monomial_coef(n,m);  
    Approx[:,i]=B*x;  
end
```

```
∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"));
```

```
∴
```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ display(Plots.plot(log.(m_values),err, label="error"))

∴ num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

∴ print("The order of convergence is estimated to be", "\n", q)

```

The order of convergence is estimated to be  
0.22324262284808913

```

∴ print("The rate of convergence is", "\n", s)

```

The rate of convergence is  
0.41982876451361223

```

∴

```

```

∴ n=32; B=monomial(n,M); Approx=zeros(M,m_size);

```

```

∴ for i=1:m_size
    m=m_values[i];
    x=monomial_coef(n,m);
    Approx[:,i]=B*x;
end

```

```

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"));

```

```

∴

```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ display(Plots.plot(log.(m_values),err, label="error"))

∴ num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

∴ print("The order of convergence is estimated to be", "\n", q)

```

The order of convergence is estimated to be  
-1.2579998394391712

```
∴ print("The rate of convergence is", "\n", s)
```

The rate of convergence is  
5.536127644117644

```
∴
```

## Problem 2.5

(Solve Problem 4, Track 1.) Instead of the monomial family of the Runge example, use the piecewise linear  $B$ -spline basis

$$A(t) = [ N_1(t) \ N_2(t) \ \dots \ N_n(t) ],$$

$$N_i(t) = \begin{cases} 0, & t < t_{i-1} \\ \frac{t - t_{i-1}}{h} & t_{i-1} \leq t < t_i \\ \frac{t_{i+1} - t}{h} & t_i \leq t < t_{i+1} \\ 0 & t_{i+1} < t \end{cases}, h = \frac{2}{m-1},$$

keeping other parameters as in Problem 1.

(The plots I got seem quite strange but I'm not sure why that would be the case.)

```
∴ function bspline(n,m,M)
    h_m = 2.0/(m-1)
    h_M = 2.0/(M-1)
    t = [(i*h_m-1) for i=0:(m+1)]
    B = zeros(M,n)
    for i=1:n
        j_min = Int(max(ceil((i-1)*h_m/h_M),1))
        j_med = Int(ceil(i*h_m/h_M)-1)
        j_max = Int(min(ceil((i+1)*h_m/h_M)-1,M))
        for j=j_min:Int(min(j_med,M))
            B[j,i] = (j*h_M-1-t[i])/h_m
        end
        if j_med+1 < M
            for j=(j_med+1):j_max
                B[j,i] = (t[i+2]-j*h_M+1)/h_m
            end
        end
    end
    return B
end
```

bspline

Julia (1.6.2) session in GNU TeXmacs

```
∴ # an alternative way to define the bspline basis
function bspline_alt(n,m,M)
    p = max(n,m)
    h_m = 2.0/(m-1)
    h_M = 2.0/(M-1)
    t = [(i*h_m-1) for i=0:(p+1)]
    T = [(i*h_M-1) for i=0:(M+1)]
    B = zeros(M,n)
    for k=1:n
        for i=1:M
            for j=1:m
                if T[i+1] >= t[k] && T[i+1] < t[k+1]
                    B[i,k]=(T[i+1]-t[k])/h_m
                elseif T[i+1] >= t[k+1] && T[i+1] < t[k+2]
                    B[i,k]=(t[k+2]-T[i+1])/h_m
                end
            end
        end
    end
    return B
end
```

bspline\_alt

∴

```
∴ function bspline_coef(n,m)
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = bspline_alt(n,m,m)
    y = runge_f.(t)
    x = A\y
    return x
end
```

bspline\_coef

∴

```
∴ for i=1:4
    local n = 2^(i+1)
    m=96
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = bspline_alt(n,m,m)
    display(Plots.plot(A, legend=false))
end
```

UndefVarError(:Plots)

```
∴ B = bspline(16,16,512); Plots.plot(B)
```

```
Plot{Plots.GRBackend() n=16}
```

```
∴
```

```
∴ M = 512; H = 2.0/(M-1); T = [(i*H-1) for i=0:(M-1)]; Y = runge_f(T);
```

```
∴ m_values = 2.^collect(4:8); m_size = length(m_values); err =  
    Array{Float64}(undef, m_size);
```

```
∴
```

```
∴ n = 4; Approx = zeros(M,m_size);
```

```
∴ for i=1:m_size  
    m = m_values[i]  
    x = bspline_coef(n,m)  
    local B = bspline(n,m,M)  
    Approx[:,i] = B*x  
end
```

```
∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))
```

```
∴ for i=1:m_size  
    err_abs = abs.(Y-Approx[:,i])  
    err[i] = dot(err_abs,err_abs)  
end
```

```
∴ err = log.(err.^0.5);
```

```
∴ display(Plots.plot(log.(m_values),err, label="error"))
```

```
∴ num=err[3:m_size] .- err[2:(m_size-1)];  
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);  
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));
```

```
∴ print("The_order_of_convergence_is_estimated_to_be","\n",q)
```

The order of convergence is estimated to be  
0.2464832341748203

```
∴ print("The_rate_of_convergence_is","\n",s)
```

The rate of convergence is  
1.1570229053972054

```
∴
```

```
∴ n = 8; Approx = zeros(M,m_size);
```

```

∴ for i=1:m_size
    m = m_values[i]
    x = bspline_coef(n,m)
    local B = bspline(n,m,M)
    Approx[:,i] = B*x
end

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ display(Plots.plot(log.(m_values),err, label="error"))

∴ num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

∴ print("The order of convergence is estimated to be", "\n", q)

The order of convergence is estimated to be
0.2464832341748203

∴ print("The rate of convergence is", "\n", s)

The rate of convergence is
1.1570229053972054

```

∴

```

∴ n = 16; Approx = zeros(M,m_size);

∴ for i=1:m_size
    m = m_values[i]
    x = bspline_coef(n,m)
    local B = bspline(n,m,M)
    Approx[:,i] = B*x
end

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ display(Plots.plot(log.(m_values),err, label="error"))

```

```

∴ num=err[3:m_size] .- err[2:(m_size-1)];
   den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
   s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

∴ print("The order of convergence is estimated to be", "\n", q)

```

The order of convergence is estimated to be  
0.2464832341748203

```

∴ print("The rate of convergence is", "\n", s)

```

The rate of convergence is  
1.1570229053972054

∴

```

∴ n = 32; Approx = zeros(M,m_size);

```

```

∴ for i=1:m_size
    m = m_values[i]
    x = bspline_coef(n,m)
    local B = bspline(n,m,M)
    Approx[:,i] = B*x
end

```

```

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

```

```

∴ err = log.(err.^0.5);

```

```

∴ display(Plots.plot(log.(m_values),err, label="error"))

```

```

∴ num=err[3:m_size] .- err[2:(m_size-1)];
   den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
   s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

```

```

∴ print("The order of convergence is estimated to be", "\n", q)

```

The order of convergence is estimated to be  
0.2464832341748203

```

∴ print("The rate of convergence is", "\n", s)

```

The rate of convergence is  
1.1570229053972054

∴

## Problem 2.6

In Problem 1, Track 1, replace the monomial basis with the Legendre polynomials, whose



samples are determined by  $QR$  decomposition  $QR = A$ . The resulting least squares problem is now

$$\min_{x \in \mathbb{R}^n} \|y - Qx\|_2.$$

```

∴ function legendre(n,m)
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:m-1]
    A = monomial(n,m)
    Q,R = qr(A)
    S = diagm(1.0 ./ Q[m,:])
    P = Q*S
    return P[:,1:n]
end

```

legendre

```

∴ function legendre_coef(n,m)
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    A = legendre(n,m)
    y = runge_f.(t)
    x = A\y
    return x
end

```

legendre\_coef

```

∴

```

```

∴ M = 512; H = 2.0/(M-1); T = [(i*H-1) for i=0:(M-1)]; Y = runge_f(T);

```

```

∴ m_values = 2 .^ collect(4:8); m_size = length(m_values); err =
    Array{Float64}(undef, m_size);

```

```

∴ for i=1:3
    local n = 2^(i+1)
    m=96
    h = 2.0/(m-1)
    t = [(i*h-1) for i=0:(m-1)]
    P = legendre(n,m)
    display(Plots.plot(P, legend=false))
end

```

```

∴

```

```

∴ n = 4; B = legendre(n,M); Approx = zeros(M,m_size);

```

```

    for i=1:m_size
        m = m_values[i]
        x = legendre_coef(n,m)
        Approx[:,i] = B*x
    end

    Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

    .

```

```

    for i=1:m_size
        err_abs = abs.(Y-Approx[:,i])
        err[i] = dot(err_abs,err_abs)
    end

    err = log.(err .^0.5);

    display(Plots.plot(log.(m_values),err, label="error"))

    num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

    print("The_order_of_convergence_is_estimated_to_be","\n",q)

```

The order of convergence is estimated to be  
0.2408324715806308

```

    print("The_rate_of_convergence_is","\n",s)

```

The rate of convergence is  
3.058988236991122

.

```

    n = 8; B = legendre(n,M); Approx = zeros(M,m_size);

    for i=1:m_size
        m = m_values[i]
        x = legendre_coef(n,m)
        Approx[:,i] = B*x
    end

    Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

    .

```

```

    for i=1:m_size
        err_abs = abs.(Y-Approx[:,i])
        err[i] = dot(err_abs,err_abs)
    end

    err = log.(err .^0.5);

```

```

∴ display(Plots.plot(log.(m_values),err, label="error"))
∴ num=err[3:m_size] .- err[2:(m_size-1)];
  den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
  s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));
∴ print("The_order_of_convergence_is_estimated_to_be","\n",q)

```

The order of convergence is estimated to be  
0.38479479155555435

```

∴ print("The_rate_of_convergence_is","\n",s)

```

The rate of convergence is  
1.521563387971189

```

∴

```

```

∴ n = 16; B = legendre(n,M); Approx = zeros(M,m_size);

```

```

∴ for i=1:m_size
    m = m_values[i]
    x = legendre_coef(n,m)
    Approx[:,i] = B*x
end

```

```

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

```

```

∴

```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

```

```

∴ err = log.(err .^0.5);

```

```

∴ display(Plots.plot(log.(m_values),err, label="error"))

```

```

∴ num=err[3:m_size] .- err[2:(m_size-1)];
  den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
  s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

```

```

∴ print("The_order_of_convergence_is_estimated_to_be","\n",q)

```

The order of convergence is estimated to be  
1.0650633160226346

```

∴ print("The_rate_of_convergence_is","\n",s)

```

The rate of convergence is  
0.6086171218702784

```

∴

```

```

∴ n = 32; B = legendre(n,M); Approx = zeros(M,m_size-1);

∴ for i=1:m_size-1
    m = m_values[i+1];
    x = legendre_coef(n,m);
    Approx[:,i] = B*x;
end

∴ Plots.plot(T,Approx); display(plot!(T,Y,label="Expected"))

∴

```

```

∴ for i=1:m_size
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ display(Plots.plot(log.(m_values),err, label="error"))

∴ num=err[3:m_size] .- err[2:(m_size-1)];
    den=err[2:m_size-1]-err[1:m_size-2]; q=sum(num ./ den)/(m_size-2);
    s=exp(sum(err[2:m_size]-q*err[1:m_size-1])/(m_size-1));

∴ print("The order of convergence is estimated to be", "\n", q)

```

The order of convergence is estimated to be  
0.2408324715806308

```

∴ print("The rate of convergence is", "\n", s)

```

The rate of convergence is  
3.058988236991122

```

∴

```

```

∴ err = zeros(m_size-1);

∴ for i=1:m_size-1
    err_abs = abs.(Y-Approx[:,i])
    err[i] = dot(err_abs,err_abs)
end

∴ err = log.(err.^0.5);

∴ a = [4; 4.5; 5];

∴ Plots.plot(log.(m_values[1:m_size-1]),err, label="error");
    display(plot!(a, (a*(-1.5) .+ 6), label="slope=-1.5"));

∴ num=err[3:m_size-1] .- err[2:(m_size-2)];
    den=err[2:m_size-2]-err[1:m_size-3]; q=sum(num ./ den)/(m_size-3);
    s=exp(sum(err[2:m_size-1]-q*err[1:m_size-2])/(m_size-2));

```

```
∴ print("The order of convergence is estimated to be", "\n", q)
```

The order of convergence is estimated to be  
1.396601993488745

```
∴ print("The rate of convergence is", "\n", s)
```

The rate of convergence is  
0.46366810375093176

```
∴
```