

Chapter 6

Approximation

Definition 6.1. Given a normed vector space Y of functions and its subspace $X \subseteq Y$. A function $\hat{\varphi} \in X$ is called the *best approximation* to $f \in Y$ from X with respect to the norm $\|\cdot\|$ iff

$$\forall \varphi \in X, \quad \|f - \hat{\varphi}\| \leq \|f - \varphi\|. \quad (6.1)$$

Example 6.1. The Chebyshev Theorem 3.32 can be restated in the format of Definition 6.1 as follows. As in Example 0.50, denote by $\mathbb{P}_n(\mathbb{R})$ the set of all polynomials with coefficients in \mathbb{R} and degree at most n . For $Y = \mathbb{P}_n(\mathbb{R})$, and $X = \mathbb{P}_{n-1}(\mathbb{R})$, the best approximation to $f(x) = -x^n$ in Y from X with respect to the max-norm $\|\cdot\|_\infty$

$$\|g\|_\infty = \max_{x \in [-1, 1]} |g(x)| \quad (6.2)$$

is $\hat{\varphi} = \frac{T_n}{2^{n-1}} - x^n$, where T_n is Chebyshev polynomial of degree n . Clearly $\hat{\varphi}$ satisfies (6.1).

Example 6.2. For $f(x) = e^x$ in $\mathcal{C}^\infty[-1, 1]$, seeking its best approximation of the form $\hat{\varphi} = \sum_{i=1}^n a_i u_i$ in the subspace $X = \text{span}\{1, x, x^2, \dots\}$ is a problem of linear approximation, where n can be any positive integer and the norm can be the max-norm (6.2), the 1-norm

$$\|g\|_1 := \int_{-1}^{+1} |g(x)| dx, \quad (6.3)$$

or the 2-norm

$$\|g\|_2 := \left(\int_{-1}^{+1} |g(x)|^2 dx \right)^{\frac{1}{2}}. \quad (6.4)$$

The three different norms are motivated differently: the max-norm corresponds to the min-max error, the 1-norm is related to the area bounded between $g(x)$ and the x -axis, and the 2-norm is related to the Euclidean distance, c.f. Section 6.4.

Example 6.3. For a simple closed curve $\gamma : [0, 1) \rightarrow \mathbb{R}^2$ and n points $\mathbf{x}_i \in \gamma$, consider a spline approximation $p : [0, 1) \rightarrow \mathbb{R}^2$ with its knots at \mathbf{x}_i 's and a scaled cumulative chordal length as in Definition 4.56. Denote by $\text{Int}(\gamma)$ as the complement of γ that always lies at the left of an observer who travels γ according to its parametrization. Then

the area difference between $\mathcal{S}_1 := \text{Int}(\gamma)$ and $\mathcal{S}_2 := \text{Int}(p)$ can be defined as

$$\|\mathcal{S}_1 \oplus \mathcal{S}_2\|_1 := \int_{\mathcal{S}_1 \oplus \mathcal{S}_2} d\mathbf{x},$$

where

$$\mathcal{S}_1 \oplus \mathcal{S}_2 := \mathcal{S}_1 \cup \mathcal{S}_2 \setminus (\mathcal{S}_1 \cap \mathcal{S}_2)$$

is the exclusive disjunction of \mathcal{S}_1 and \mathcal{S}_2 .

The minimization of this area difference can be formulated by a best approximation problem based on the 1-norm.

Theorem 6.2. Suppose X is a finite-dimensional subspace of a normed space $(Y, \|\cdot\|)$. Then we have

$$\forall y \in Y, \exists \hat{\varphi} \in X \text{ s.t. } \forall \varphi \in X, \|\hat{\varphi} - y\| \leq \|\varphi - y\|. \quad (6.5)$$

Proof. For a given $y \in Y$, define a closed ball

$$B_y := \{x \in X : \|x\| \leq 2\|y\|\}.$$

Clearly $0 \in B_y$, and the distance from y to B_y is

$$\text{dist}(y, B_y) := \inf_{x \in B_y} \|y - x\| \leq \|y - 0\| = \|y\|.$$

By definition, any $z \in X$, $z \notin B_y$ must satisfy $\|z\| > 2\|y\|$, and thus

$$\|z - y\| \geq \|z\| - \|y\| > \|y\|.$$

Therefore, if a best approximation to y exists, it must be in B_y . As a subspace of X , B_y is finite dimensional, closed, and bounded, hence B_y is compact. The extreme value theorem states that a continuous scalar function attains its minimum and maximum on a compact set. A norm is a continuous function, hence the function $d : B_y \rightarrow \mathbb{R}^+ \cup \{0\}$ given by $d(x) = \|x - y\|$ must attain its minimum on B_y . \square

Definition 6.3 (L^p functions). Let $p > 0$. The class of functions $f(x)$ which are measurable and for which $|f(x)|^p$ is Lebesgue integrable over $[a, b]$ is known as $L^p[a, b]$. If $p = 1$, the class is denoted by $L[a, b]$.

Theorem 6.4. For a weight function $\rho(x) \in L[a, b]$, define

$$L_\rho^2[a, b] := \{f(x) \in L[a, b] : \rho(x)|f(x)|^2 \in L[a, b]\}. \quad (6.6)$$

Then $L_\rho^2[a, b]$ is a vector space. If we further require that $\forall x \in (a, b), \rho(x) > 0$, then the vector space $L_\rho^2[a, b]$ with

$$\langle u, v \rangle = \int_a^b \rho(t) u(t) \overline{v(t)} dt \quad (6.7)$$

is an inner product space over \mathbb{R} ; the set $L^2_\rho[a, b]$ with

$$\|u\|_2 = \left(\int_a^b \rho(t) |u(t)|^2 dt \right)^{\frac{1}{2}} \quad (6.8)$$

is a normed vector space over \mathbb{R} .

Proof. This follows from Definitions 0.69, 0.87, and 0.89. \square

Definition 6.5. The *least-square approximation* on $L^2_\rho[a, b]$ is a best approximation problem with the norm in (6.1) set to that in (6.8).

6.1 Orthonormal systems

Definition 6.6. A subset S of an inner product space X is called *orthonormal* if

$$\forall u, v \in S, \quad \langle u, v \rangle = \begin{cases} 0 & \text{if } u \neq v, \\ 1 & \text{if } u = v. \end{cases} \quad (6.9)$$

Example 6.4. The standard basis vectors in \mathbb{R}^n are orthonormal.

Example 6.5. The Chebyshev polynomials of the first kind as in Definition 3.28 are orthogonal with respect to (6.7) where $a = -1, b = 1, \rho = \frac{1}{\sqrt{1-x^2}}$. However, they do not satisfy the second case in (6.9).

Theorem 6.7. Any finite set of nonzero orthogonal elements u_1, u_2, \dots, u_n is linearly independent.

Proof. This is easily proven by contradiction using Definitions 0.76 and 6.6. \square

Definition 6.8. The *Gram-Schmidt process* takes in a finite or infinite independent list (u_1, u_2, \dots) and output two other lists (v_1, v_2, \dots) and (u_1^*, u_2^*, \dots) by

$$v_{n+1} = u_{n+1} - \sum_{k=1}^n \langle u_{n+1}, u_k^* \rangle u_k^*, \quad (6.10a)$$

$$u_{n+1}^* = v_{n+1} / \|v_{n+1}\|, \quad (6.10b)$$

with the recursion basis as $v_1 = u_1, u_1^* = v_1 / \|v_1\|$.

Theorem 6.9. For a finite or infinite independent list (u_1, u_2, \dots) , the Gram-Schmidt process yields constants

$$\begin{matrix} a_{11} \\ a_{21} & a_{22} \\ a_{31} & a_{32} & a_{33} \\ \vdots \end{matrix}$$

such that $a_{kk} = \frac{1}{\|v_k\|} > 0$ and the elements u_1^*, u_2^*, \dots

$$\begin{matrix} u_1^* = a_{11}u_1 \\ u_2^* = a_{21}u_1 + a_{22}u_2 \\ u_3^* = a_{31}u_1 + a_{32}u_2 + a_{33}u_3 \\ \vdots \end{matrix} \quad (6.11)$$

are orthonormal.

Proof. By Definition 6.8, the formulae (6.10) can be rewritten in the form of (6.11). It is clear from (6.10b) that u_{n+1}^* is normal. We show u_{n+1}^* is orthogonal to $u_n^*, u_{n-1}^*, \dots, u_1^*$ by induction. The induction base holds because

$$\begin{aligned} \langle v_2, u_1^* \rangle &= \langle u_2 - \langle u_2, u_1^* \rangle u_1^*, u_1^* \rangle \\ &= \langle u_2, u_1^* \rangle - \langle u_2, u_1^* \rangle \langle u_1^*, u_1^* \rangle = 0, \end{aligned}$$

where the second step follows from (IP-3) in Definition 0.87 and the third step from u_1^* being normal. The inductive step also holds because for any $j < n + 1$ we have

$$\begin{aligned} \langle v_{n+1}, u_j^* \rangle &= \left\langle u_{n+1} - \sum_{k=1}^n \langle u_{n+1}, u_k^* \rangle u_k^*, u_j^* \right\rangle \\ &= \langle u_{n+1}, u_j^* \rangle - \sum_{k=1}^n \langle u_{n+1}, u_k^* \rangle \langle u_k^*, u_j^* \rangle \\ &= \langle u_{n+1}, u_j^* \rangle - \langle u_{n+1}, u_j^* \rangle = 0, \end{aligned}$$

where the third step follows from the induction hypothesis, i.e., $\langle u_k^*, u_j^* \rangle$ is 1 if $k = j$ and 0 otherwise. It remains to show $a_{kk} = \frac{1}{\|v_k\|}$, which holds because

$$\begin{aligned} 1 = \langle u_n^*, u_n^* \rangle &= \langle a_{nn}u_n, u_n^* \rangle + \left\langle \sum_{i=1}^{n-1} a_{ni}u_i, u_n^* \right\rangle \\ &= a_{nn} \langle u_n, u_n^* \rangle = a_{nn} \left\langle u_n, \frac{v_n}{\|v_n\|} \right\rangle = a_{nn} \|v_n\|, \end{aligned}$$

where the second step follows from the n th equation of (6.11), the third step from (6.10a) and the conclusion just proved, the fourth step from (6.10b), and the last step from Definitions 0.87 and 0.89. \square

Corollary 6.10. For a finite or infinite independent list (u_1, u_2, \dots) , we can find constants

$$\begin{matrix} b_{11} \\ b_{21} & b_{22} \\ b_{31} & b_{32} & b_{33} \\ \vdots \end{matrix}$$

and an orthonormal list (u_1^*, u_2^*, \dots) such that $b_{ii} > 0$ and

$$\begin{aligned} u_1 &= b_{11}u_1^* \\ u_2 &= b_{21}u_1^* + b_{22}u_2^* \\ u_3 &= b_{31}u_1^* + b_{32}u_2^* + b_{33}u_3^* \\ &\vdots \end{aligned} \quad (6.12)$$

Proof. This follows from (6.11) and that a lower-triangular matrix with positive diagonal elements is invertible. \square

Corollary 6.11. In Theorem 6.9, we have $\langle u_n^*, u_i \rangle = 0$ for each $i = 1, 2, \dots, n - 1$.

Proof. By Corollary 6.10, each u_i can be expressed as

$$u_i = \sum_{k=1}^i b_{ik}u_k^*.$$

Inner product the above equation with u_n^* , apply the orthogonal conditions, and we reach the conclusion. \square

Definition 6.12. Using the Gram-Schmidt orthonormalizing process with the inner product (6.7), we obtain from the independent list of monomials $(1, x, x^2, \dots)$ the following *classic orthonormal polynomials*:

	a	b	$\rho(x)$
Chebyshev polynomials of the first kind	-1	1	$\frac{1}{\sqrt{1-x^2}}$
Chebyshev polynomials of the second kind	-1	1	$\sqrt{1-x^2}$
Legendre polynomials	-1	1	1
Jacobi polynomials	-1	1	$(1-x)^\alpha(1+x)^\beta$
Laguerre polynomials	0	$+\infty$	$x^\alpha e^{-x}$
Hermite polynomials	$-\infty$	$+\infty$	e^{-x^2}

where $\alpha, \beta > -1$ for Jacobi polynomials and $\alpha > -1$ for Laguerre polynomials.

Example 6.6. We compute the first 3 Legendre polynomials using the Gram-Schmidt process.

$$\begin{aligned}
 u_1 &= 1, \quad v_1 = 1, \quad \|v_1\|^2 = \int_{-1}^{+1} dx = 2, \quad u_1^* = \frac{1}{\sqrt{2}}. \\
 u_2 &= x, \quad v_2 = x - \left\langle x, \frac{1}{\sqrt{2}} \right\rangle \frac{1}{\sqrt{2}} = x, \quad \|v_2\|^2 = \frac{2}{3}, \\
 u_2^* &= \sqrt{\frac{3}{2}}x. \\
 v_3 &= x^2 - \left\langle x^2, \sqrt{\frac{3}{2}}x \right\rangle \sqrt{\frac{3}{2}}x - \left\langle x^2, \frac{1}{\sqrt{2}} \right\rangle \frac{1}{\sqrt{2}} = x^2 - \frac{1}{3}, \\
 \|v_3\|^2 &= \int_{-1}^{+1} \left(x^2 - \frac{1}{3} \right)^2 dx = \frac{8}{45}, \\
 u_3^* &= \frac{3}{4}\sqrt{10} \left(x^2 - \frac{1}{3} \right).
 \end{aligned}$$

6.2 Fourier expansions

Definition 6.13. Let (u_1^*, u_2^*, \dots) be a finite or infinite orthonormal list. The *orthogonal expansion* or *Fourier expansion* for an arbitrary w is the series

$$w \sim \sum_n \langle w, u_n^* \rangle u_n^*, \quad (6.13)$$

where the constants $\langle w, u_n^* \rangle$ are known as the *Fourier coefficients* of w and the term $\langle w, u_n^* \rangle u_n^*$ the *projection* of w on u_n^* . The *error of the Fourier expansion* of w with respect to (u_1^*, u_2^*, \dots) is simply $\sum_n \langle w, u_n^* \rangle u_n^* - w$.

Example 6.7. With the Euclidean inner product in Definition 0.88, we select orthonormal vectors in \mathbb{R}^3 as

$$u_1^* = (1, 0, 0)^T, \quad u_2^* = (0, 1, 0)^T, \quad u_3^* = (0, 0, 1)^T.$$

For the vector $w = (a, b, c)^T$, the Fourier coefficients are

$$\langle w, u_1^* \rangle = a, \quad \langle w, u_2^* \rangle = b, \quad \langle w, u_3^* \rangle = c,$$

and the projections of w onto u_1^* and u_2^* are

$$\langle w, u_1^* \rangle u_1^* = (a, 0, 0)^T, \quad \langle w, u_2^* \rangle u_2^* = (0, b, 0)^T.$$

The Fourier expansion of w is

$$w = \langle w, u_1^* \rangle u_1^* + \langle w, u_2^* \rangle u_2^* + \langle w, u_3^* \rangle u_3^*,$$

with the error of Fourier expansion as 0; see Theorem 6.14.

Exercise 6.8. With the following orthonormal list in $L_{\rho=1}^2[-\pi, \pi]$,

$$\frac{1}{\sqrt{2\pi}}, \frac{\sin x}{\sqrt{\pi}}, \frac{\cos x}{\sqrt{\pi}}, \dots, \frac{\sin(nx)}{\sqrt{\pi}}, \frac{\cos(nx)}{\sqrt{\pi}}, \dots, \quad (6.14)$$

derive the *Fourier series* of a function $f(x)$ as

$$f(x) \sim \frac{a_0}{2} + \sum_{k=1}^{+\infty} (a_k \cos kx + b_k \sin kx), \quad (6.15)$$

where the coefficients are

$$a_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos kx dx, \quad b_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin kx dx.$$

Theorem 6.14. Let u_1, u_2, \dots, u_n be linearly independent and let u_i^* be the u_i 's orthonormalized by the Gram-Schmidt process. If $w = \sum_{i=1}^n a_i u_i$, then

$$w = \sum_{i=1}^n \langle w, u_i^* \rangle u_i^*, \quad (6.16)$$

i.e. w is equal to its Fourier expansion.

Proof. By the condition $w = \sum_{i=1}^n a_i u_i$ and Corollary 6.10, we can express w as a linear combination of u_i^* 's,

$$w = \sum_{i=1}^n c_i u_i^*.$$

Then the orthogonality of u_i^* 's implies

$$\forall k = 1, 2, \dots, n, \quad \langle u_k^*, w \rangle = c_k,$$

which completes the proof. \square

Theorem 6.15 (Minimum properties of Fourier expansions). Let u_1^*, u_2^*, \dots be an orthonormal system and let w be arbitrary. Then

$$\left\| w - \sum_{i=1}^N \langle w, u_i^* \rangle u_i^* \right\| \leq \left\| w - \sum_{i=1}^N a_i u_i^* \right\|, \quad (6.17)$$

for any selection of constants a_1, a_2, \dots, a_N .

Proof. With the shorthand notation $\sum_i = \sum_{i=1}^N$, we deduce

from Definition 0.87 and properties of inner products

$$\begin{aligned}
\left\| w - \sum_i a_i u_i^* \right\|^2 &= \left\langle w - \sum_i a_i u_i^*, w - \sum_i a_i u_i^* \right\rangle \\
&= \langle w, w \rangle - \left\langle w, \sum_i a_i u_i^* \right\rangle - \left\langle \sum_i a_i u_i^*, w \right\rangle \\
&\quad + \left\langle \sum_i a_i u_i^*, \sum_i a_i u_i^* \right\rangle \\
&= \langle w, w \rangle - \sum_i \overline{a_i} \langle w, u_i^* \rangle - \sum_i a_i \langle u_i^*, w \rangle \\
&\quad + \sum_i \sum_j a_i \overline{a_j} \langle u_i^*, u_j^* \rangle \\
&= \langle w, w \rangle - \sum_i \overline{a_i} \langle w, u_i^* \rangle - \sum_i a_i \langle u_i^*, w \rangle + \sum_i |a_i|^2 \\
&\quad - \sum_i \langle u_i^*, w \rangle \langle w, u_i^* \rangle + \sum_i \langle u_i^*, w \rangle \langle w, u_i^* \rangle \\
&= \|w\|^2 - \sum_i |\langle w, u_i^* \rangle|^2 + \sum_i |a_i - \langle w, u_i^* \rangle|^2, \quad (6.18)
\end{aligned}$$

where “ $|\cdot|$ ” denotes the modulus of a complex number. The first two terms are independent of a_i . Therefore $\|w - \sum_i a_i u_i^*\|^2$ is minimized only when $a_i = \langle w, u_i^* \rangle$. \square

Corollary 6.16. Let (u_1, u_2, \dots, u_n) be an independent list. The fundamental problem of linearly approximating an arbitrary vector w is solved by the best approximation $\hat{\varphi} = \sum_k \langle w, u_k^* \rangle u_k^*$ where u_k^* 's are the u_k 's orthonormalized by the Gram-Schmidt process. The error norm is

$$\|w - \hat{\varphi}\|^2 := \min_{a_k} \left\| w - \sum_{k=1}^n a_k u_k \right\|^2 = \|w\|^2 - \sum_{k=1}^n |\langle w, u_k^* \rangle|^2. \quad (6.19)$$

Proof. This follows directly from (6.18). \square

Corollary 6.17 (Bessel inequality). If $u_1^*, u_2^*, \dots, u_N^*$ are orthonormal, then, for an arbitrary w ,

$$\sum_{i=1}^N |\langle w, u_i^* \rangle|^2 \leq \|w\|^2. \quad (6.20)$$

Proof. This follows directly from Corollary 6.16 and the real positivity of a norm. \square

Corollary 6.18. The Gram-Schmidt process in Definition 6.8 satisfies

$$\forall n \in \mathbb{N}^+, \quad \|v_{n+1}\|^2 = \|u_{n+1}\|^2 - \sum_{k=1}^n |\langle u_{n+1}, u_k^* \rangle|^2. \quad (6.21)$$

Proof. By (6.10a), each v_{n+1} can be regarded as the error of Fourier expansion of u_{n+1} with respect to the orthonormal list $(u_1^*, u_2^*, \dots, u_n^*)$. In Corollary 6.16, identifying w with u_{n+1} completes the proof. \square

Example 6.9. Consider the problem in Example 6.2 in the sense of least square approximation with the weight function $\rho = 1$. It is equivalent to

$$\min_{a_i} \int_{-1}^{+1} \left(e^x - \sum_{i=0}^n a_i x^i \right)^2 dx. \quad (6.22)$$

For $n = 1, 2$, use the Legendre polynomials derived in Example 6.6:

$$u_1^* = \frac{1}{\sqrt{2}}, \quad u_2^* = \sqrt{\frac{3}{2}}x, \quad u_3^* = \frac{1}{4}\sqrt{10}(3x^2 - 1),$$

and we have the Fourier coefficients of e^x as

$$\begin{aligned}
b_0 &= \int_{-1}^{+1} \frac{1}{\sqrt{2}} e^x dx = \frac{1}{\sqrt{2}} \left(e - \frac{1}{e} \right), \\
b_1 &= \int_{-1}^{+1} \sqrt{\frac{3}{2}} x e^x dx = \sqrt{6} e^{-1}, \\
b_2 &= \int_{-1}^{+1} \frac{1}{4} \sqrt{10} (3x^2 - 1) e^x dx = \frac{\sqrt{10}}{2} \left(e - \frac{7}{e} \right).
\end{aligned}$$

The minimizing polynomials are thus

$$\hat{\varphi}_n = \begin{cases} \frac{1}{2e}(e^2 - 1) + \frac{3}{e}x & n = 1; \\ \hat{\varphi}_1 + \frac{5}{4e}(e^2 - 7)(3x^2 - 1) & n = 2. \end{cases} \quad (6.23)$$

6.3 The normal equations

Theorem 6.19. Let $u_1, u_2, \dots, u_n \in X$ be linearly independent and let u_i^* be the u_i 's orthonormalized by the Gram-Schmidt process. Then, for any element w ,

$$\forall j = 1, 2, \dots, n, \quad \left(w - \sum_{k=1}^n \langle w, u_k^* \rangle u_k^* \right) \perp u_j^*, \quad (6.24)$$

where “ \perp ” denotes orthogonality.

Proof. Take the inner product of the two vectors and apply the conditions on orthonormal systems. \square

Corollary 6.20. Let $u_1, u_2, \dots, u_n \in X$ be linearly independent. If $\hat{\varphi} = \sum_{k=1}^n a_k u_k$ is the best linear approximant to w , then

$$\forall j = 1, 2, \dots, n, \quad (w - \hat{\varphi}) \perp u_j. \quad (6.25)$$

Proof. Since $\hat{\varphi} = \sum_{k=1}^n a_k u_k$ is the best linear approximant to w , Theorem 6.15 implies that

$$\sum_{k=1}^n a_k u_k = \sum_{k=1}^n \langle w, u_k^* \rangle u_k^*.$$

Corollary 6.10 and Theorem 6.19 complete the proof. \square

Definition 6.21. Let u_1, u_2, \dots, u_n be a sequence of elements in an inner product space. The $n \times n$ matrix

$$\begin{aligned}
G &= G(u_1, u_2, \dots, u_n) = (\langle u_i, u_j \rangle) \\
&= \begin{bmatrix} \langle u_1, u_1 \rangle & \langle u_1, u_2 \rangle & \dots & \langle u_1, u_n \rangle \\ \langle u_2, u_1 \rangle & \langle u_2, u_2 \rangle & \dots & \langle u_2, u_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_n, u_1 \rangle & \langle u_n, u_2 \rangle & \dots & \langle u_n, u_n \rangle \end{bmatrix} \quad (6.26)
\end{aligned}$$

is the *Gram matrix* of u_1, u_2, \dots, u_n . Its determinant

$$g = g(u_1, u_2, \dots, u_n) = \det(\langle u_i, u_j \rangle) \quad (6.27)$$

is the *Gram determinant*.

Lemma 6.22. Let $w_i = \sum_{j=1}^n a_{ij} u_j$ for $i = 1, 2, \dots, n$. Let $A = (a_{ij})$ and its conjugate transpose $A^H = (\overline{a_{ji}})$. Then we have

$$G(w_1, w_2, \dots, w_n) = AG(u_1, u_2, \dots, u_n)A^H \quad (6.28)$$

and

$$g(w_1, w_2, \dots, w_n) = |\det A|^2 g(u_1, u_2, \dots, u_n). \quad (6.29)$$

Proof. The inner product of u_i and w_j yields

$$\begin{aligned} & \begin{bmatrix} \langle u_1, w_1 \rangle & \langle u_1, w_2 \rangle & \dots & \langle u_1, w_n \rangle \\ \langle u_2, w_1 \rangle & \langle u_2, w_2 \rangle & \dots & \langle u_2, w_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_n, w_1 \rangle & \langle u_n, w_2 \rangle & \dots & \langle u_n, w_n \rangle \end{bmatrix} \\ &= \begin{bmatrix} \langle u_1, u_1 \rangle & \langle u_1, u_2 \rangle & \dots & \langle u_1, u_n \rangle \\ \langle u_2, u_1 \rangle & \langle u_2, u_2 \rangle & \dots & \langle u_2, u_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_n, u_1 \rangle & \langle u_n, u_2 \rangle & \dots & \langle u_n, u_n \rangle \end{bmatrix} \begin{bmatrix} \overline{a_{11}} & \dots & \overline{a_{n1}} \\ \overline{a_{12}} & \dots & \overline{a_{n2}} \\ \vdots & \ddots & \vdots \\ \overline{a_{1n}} & \dots & \overline{a_{nn}} \end{bmatrix} \\ &= G(u_1, u_2, \dots, u_n) A^H. \end{aligned}$$

Therefore (6.28) holds since

$$\begin{aligned} G(w_1, w_2, \dots, w_n) &= \begin{bmatrix} \langle w_1, w_1 \rangle & \langle w_1, w_2 \rangle & \dots & \langle w_1, w_n \rangle \\ \langle w_2, w_1 \rangle & \langle w_2, w_2 \rangle & \dots & \langle w_2, w_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle w_n, w_1 \rangle & \langle w_n, w_2 \rangle & \dots & \langle w_n, w_n \rangle \end{bmatrix} \\ &= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ a_{21} & \dots & a_{2n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} \langle u_1, w_1 \rangle & \langle u_1, w_2 \rangle & \dots & \langle u_1, w_n \rangle \\ \langle u_2, w_1 \rangle & \langle u_2, w_2 \rangle & \dots & \langle u_2, w_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_n, w_1 \rangle & \langle u_n, w_2 \rangle & \dots & \langle u_n, w_n \rangle \end{bmatrix} \\ &= AG(u_1, u_2, \dots, u_n) A^H. \end{aligned}$$

The following properties of complex conjugate are well known:

$$\overline{z + w} = \overline{z} + \overline{w}, \quad \overline{zw} = \overline{z} \overline{w}.$$

Then the identity $\det(A) = \det(A^T)$ and the Leibniz formula of determinants (0.86) yields

$$\overline{\det A} = \det A^T = \sum_{\sigma \in S_n} \operatorname{sgn}(\sigma) \prod_{i=1}^n \overline{a_{\sigma(i), i}} = \det A^H.$$

Take the determinant of (6.28), apply the identity $\det(AB) = \det(A) \det(B)$, and we have (6.29). \square

Theorem 6.23. For nonzero elements $u_1, u_2, \dots, u_n \in X$, we have

$$0 \leq g(u_1, u_2, \dots, u_n) \leq \prod_{k=1}^n \|u_k\|^2, \quad (6.30)$$

where the lower equality holds if and only if u_1, u_2, \dots, u_n are linearly dependent and the upper equality holds if and only if they are orthogonal.

Proof. Suppose u_1, u_2, \dots, u_n are linearly dependent. Then we can find constants c_1, c_2, \dots, c_n such that $\sum_{i=1}^n c_i u_i = \mathbf{0}$ with at least one constant c_j being nonzero. Construct vectors

$$w_k = \begin{cases} \sum_{i=1}^n c_i u_i = \mathbf{0}, & k = j; \\ u_k, & k \neq j. \end{cases}$$

We have $g(w_1, w_2, \dots, w_n) = 0$ because $\langle w_j, w_k \rangle = 0$ for each k . By the Laplace theorem, we expand the determinant of $C = (c_{ij})$ according to minors of its j th row:

$$\begin{aligned} \det(C) &= \det \begin{bmatrix} 1 & 0 & \dots & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ c_1 & c_2 & \dots & c_j & \dots & c_n \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \\ &= 0 + \dots + 0 + c_j + 0 + \dots + 0 = c_j \neq 0, \end{aligned}$$

where the determinant of each minor matrix M_i of c_i with $i \neq j$ is zero because each M_i has a row of all zeros. Then Lemma 6.22 yields $g(u_1, u_2, \dots, u_n) = 0$.

Now suppose u_1, u_2, \dots, u_n are linearly independent. Theorem 6.9 yields constants a_{ij} such that $a_{kk} > 0$ and the following vectors are orthonormal:

$$u_k^* = \sum_{i=1}^k a_{ki} u_i.$$

Then Definition 6.21 implies $g(u_1^*, u_2^*, \dots, u_n^*) = 1$. Also, we have $\det(a_{ij}) = \prod_{k=1}^n a_{kk}$ because the matrix (a_{ij}) is triangular. It then follows from Lemma 6.22 that

$$g(u_1, u_2, \dots, u_n) = \prod_{k=1}^n \frac{1}{a_{kk}^2} > 0. \quad (6.31)$$

Since the list of vectors (u_1, u_2, \dots, u_n) is either dependent or independent, the arguments so far show that $g(u_1, u_2, \dots, u_n) = 0$ if and only if u_1, u_2, \dots, u_n are linearly dependent.

Suppose u_1, u_2, \dots, u_n are orthogonal. By Definition 6.21, $G(u_1, u_2, \dots, u_n)$ is a diagonal matrix with $\|u_k\|^2$ on the diagonals. Hence the orthogonality of u_k 's implies

$$g(u_1, u_2, \dots, u_n) = \prod_{k=1}^n \|u_k\|^2. \quad (6.32)$$

For the converse statement, suppose (6.32) holds. Then u_1, u_2, \dots, u_n must be independent because otherwise it would contradict the lower equality proved as above. Apply the Gram-Schmidt process to (u_1, u_2, \dots, u_n) and we know from Theorem 6.9 that $\frac{1}{a_{kk}} = \|v_k\|$. Set the length of the list in Theorem 6.9 to $1, 2, \dots, n$ and we know from (6.31) and (6.32) that

$$\forall k = 1, 2, \dots, n, \quad \|u_k\|^2 = \|v_k\|^2. \quad (6.33)$$

Then Corollary 6.18 and (6.33) imply

$$\forall k = 1, 2, \dots, n, \quad \sum_{j=1}^{k-1} |\langle u_k, u_j^* \rangle|^2 = 0,$$

which further implies

$$\forall k = 1, 2, \dots, n, \quad \forall j = 1, 2, \dots, k-1, \quad \langle u_k, u_j^* \rangle = 0,$$

which, together with Corollary 6.10, implies the orthogonality of u_k 's. Finally, we remark that the maximum of $g(u_1, u_2, \dots, u_n)$ is indeed $\prod_{k=1}^n \|u_k\|^2$ because of (6.31), $\frac{1}{a_{kk}} = \|v_k\|$, and Corollary 6.18. \square

Theorem 6.24. Let $\hat{\varphi} = \sum_{i=1}^n a_i u_i$ be the best approximation to w constructed from the list of independent vectors (u_1, u_2, \dots, u_n) . Then the coefficients

$$\mathbf{a} = [a_1, a_2, \dots, a_n]^T$$

are uniquely determined from the linear system of *normal equations*,

$$G(u_1, u_2, \dots, u_n)^T \mathbf{a} = \mathbf{c}, \quad (6.34)$$

where $\mathbf{c} = [\langle w, u_1 \rangle, \langle w, u_2 \rangle, \dots, \langle w, u_n \rangle]^T$.

Proof. Corollary 6.20 yields

$$\langle w, u_j \rangle = \sum_{k=1}^n a_k \langle u_k, u_j \rangle,$$

which is simply the j th equation of (6.34). The uniqueness of the coefficients follows from Theorem 6.23 and Cramer's rule. \square

Example 6.10. Solve Example 6.9 by normal equations.

To find the best approximation $\hat{\varphi} = a_0 + a_1 x + a_2 x^2$ to e^x from the linearly independent list $(1, x, x^2)$, we first construct the Gram matrix from (6.26), (6.7), and $\rho = 1$:

$$G(1, x, x^2) = \begin{bmatrix} \langle 1, 1 \rangle & \langle 1, x \rangle & \langle 1, x^2 \rangle \\ \langle x, 1 \rangle & \langle x, x \rangle & \langle x, x^2 \rangle \\ \langle x^2, 1 \rangle & \langle x^2, x \rangle & \langle x^2, x^2 \rangle \end{bmatrix} = \begin{bmatrix} 2 & 0 & \frac{2}{3} \\ 0 & \frac{2}{3} & 0 \\ \frac{2}{3} & 0 & \frac{2}{5} \end{bmatrix}.$$

We then calculate the vector

$$\mathbf{c} = \begin{bmatrix} \langle e^x, 1 \rangle \\ \langle e^x, x \rangle \\ \langle e^x, x^2 \rangle \end{bmatrix} = \begin{bmatrix} e - 1/e \\ 2/e \\ e - 5/e \end{bmatrix}.$$

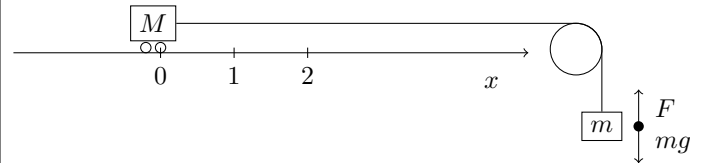
The normal equations then yields

$$a_0 = \frac{3(11 - e^2)}{4e}, \quad a_1 = \frac{3}{e}, \quad a_2 = \frac{15(e^2 - 7)}{4e}.$$

With these values, it is easily verified that the best approximation $\hat{\varphi} = a_0 + a_1 x + a_2 x^2$ equals that in (6.23).

6.4 Discrete least squares (DLS)

Example 6.11 (An experiment on Newton's second law by discrete least squares). A cart with mass M is pulled along a horizontal track by a cable attached to a weight of mass m_j through a pulley.



Neglecting the friction of the track and the pulley system, we have from Newton's second law

$$m_j g = (m_j + M) a = (m_j + M) \frac{d^2 x}{dt^2}.$$

A series of experiments can be designed to test the hypothesis of Newton's second law.

- (i) For fixed M and m_j , we measure a number of data points (t_i, x_i) by recording the position of the cart with a high-speed camera.
- (ii) Fit a quadratic polynomial $p(t) = c_0 + c_1 t + c_2 t^2$ by minimizing the total length squared,

$$\min \sum_i (x_i - p(t_i))^2.$$

- (iii) Take $a_j = 2c_2$ as the experimental result of acceleration for the force $F_j = m_j(g - a_j)$.
- (iv) Change the weight m_j and repeat steps (i)-(iii) a number of times to get data points (a_j, F_j) .
- (v) Fit a linear polynomial $f(x) = c_0 + c_1 x$ by minimizing the total length squared,

$$\min \sum_j (F_j - f(a_j))^2.$$

One verifies Newton's second law by showing that the data fitting result c_1 is very close to M . Note that the expressions in steps (ii) and (v) justify the name "least squares."

6.4.1 Reusing the formalism

Definition 6.25. Define a function $\lambda : \mathbb{R} \rightarrow \mathbb{R}$

$$\lambda(t) = \begin{cases} 0 & \text{if } t \in (-\infty, a), \\ \int_a^t \rho(\tau) d\tau & \text{if } t \in [a, b], \\ \int_a^b \rho(\tau) d\tau & \text{if } t \in (b, +\infty). \end{cases} \quad (6.35)$$

Then a corresponding *continuous measure* $d\lambda$ can be defined as

$$d\lambda = \begin{cases} \rho(t) dt & \text{if } t \in [a, b], \\ 0 & \text{otherwise,} \end{cases} \quad (6.36)$$

where the *support of the continuous measure* $d\lambda$ is the interval $[a, b]$.

Definition 6.26. The *discrete measure* or the *Dirac measure* associated with the point set $\{t_1, t_2, \dots, t_N\}$ is a measure $d\lambda$ that is nonzero only at the points t_i and has the value ρ_i there. The *support of the discrete measure* is the set $\{t_1, t_2, \dots, t_N\}$.

Definition 6.27. The *Heaviside function* is the truncated power function with exponent 0,

$$H(x) = x_+^0 = \begin{cases} 1 & \text{if } x \geq 0, \\ 0 & \text{if } x < 0. \end{cases} \quad (6.37)$$

Lemma 6.28. For a function $u : \mathbb{R} \rightarrow \mathbb{R}$, define

$$\lambda(t) = \sum_{i=1}^N \rho_i H(t - t_i), \quad (6.38)$$

and we have

$$\int_{\mathbb{R}} u(t) d\lambda = \sum_{i=1}^N \rho_i u(t_i). \quad (6.39)$$

Proof. The *Dirac Delta function*, $\delta(x)$, is roughly a generalized function that satisfies

$$\delta(x) = \begin{cases} +\infty & x = 0, \\ 0 & x \neq 0. \end{cases} \quad (6.40)$$

Note: the above definition of $\delta(x)$ is heuristic. A rigorous one should employ the concept of measures.

Useful properties of $\delta(x)$ include

$$\int_{-\infty}^{+\infty} \delta(x) dx = 1, \quad (6.41)$$

$$\int_0^x \delta(t) dt = H(x), \quad (6.42)$$

$$\int_{-\infty}^{+\infty} f(t) \delta(t - t_0) dt = f(t_0). \quad (6.43)$$

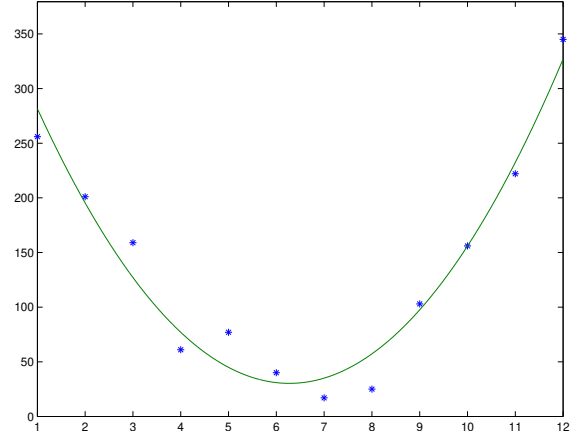
Then (6.38), (6.42), and (6.43) yield

$$\int_{\mathbb{R}} u(t) d\lambda = \int_{\mathbb{R}} \sum_{i=1}^N \rho_i \delta(t - t_i) u(t) dt = \sum_{i=1}^N \rho_i u(t_i). \quad \square$$

6.4.2 DLS via normal equations

Example 6.12. Consider a table of sales record.

x	1	2	3	4	5	6
y	256	201	159	61	77	40
x	7	8	9	10	11	12
y	17	25	103	156	222	345



From the plot of the discrete data, it appears that a quadratic polynomial would be a good fit. Hence we formulate the least square problem as finding the coefficients of a quadratic polynomial to minimize the following error,

$$\sum_{i=1}^{12} \left(y_i - \sum_{j=0}^2 a_j x_i^j \right)^2.$$

Reusing the procedures in Example 6.10, we have

$$G(1, x, x^2) = \begin{bmatrix} \langle 1, 1 \rangle & \langle 1, x \rangle & \langle 1, x^2 \rangle \\ \langle x, 1 \rangle & \langle x, x \rangle & \langle x, x^2 \rangle \\ \langle x^2, 1 \rangle & \langle x^2, x \rangle & \langle x^2, x^2 \rangle \end{bmatrix} = \begin{bmatrix} 12 & 78 & 650 \\ 78 & 650 & 6084 \\ 650 & 6084 & 60710 \end{bmatrix},$$

$$\mathbf{c} = \begin{bmatrix} \langle y, 1 \rangle \\ \langle y, x \rangle \\ \langle y, x^2 \rangle \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{12} y_i \\ \sum_{i=1}^{12} y_i x_i \\ \sum_{i=1}^{12} y_i x_i^2 \end{bmatrix} = \begin{bmatrix} 1662 \\ 11392 \\ 109750 \end{bmatrix}.$$

Then the normal equations yield

$$\mathbf{a} = G^{-1} \mathbf{c} = [386.00, -113.43, 9.04]^T.$$

The corresponding polynomial is plotted in the figure.

6.4.3 DLS via QR decomposition

Definition 6.29. A matrix $A \in \mathbb{R}^{n \times n}$ is *orthogonal* iff $A^T A = I$.

Definition 6.30. A matrix A is *upper triangular* iff

$$\forall i, j, \quad i > j \Rightarrow a_{i,j} = 0.$$

Similarly, a matrix A is *lower triangular* iff

$$\forall i, j, \quad i < j \Rightarrow a_{i,j} = 0.$$

Theorem 6.31 (QR factorization). For any $A \in \mathbb{R}^{m \times n}$, there exists an orthogonal matrix $Q \in \mathbb{R}^{m \times m}$ and an upper triangular matrix $R \in \mathbb{R}^{m \times n}$ so that $A = QR$.

Proof. Rewrite $A = [\xi_1, \xi_2, \dots, \xi_n] \in \mathbb{R}^{m \times n}$ and denote by r the column rank of A . Construct a rank- r matrix

$$A_r = [u_1, u_2, \dots, u_r]$$

by the following steps.

(S-1) Set $u_1 = \xi_{k_1}$ where k_1 satisfies $\forall \ell < k_1, \xi_\ell = \mathbf{0}$.

(S-2) For each $j = 2, \dots, r$, set $u_j = \xi_{k_j}$ where k_j satisfies that $K_j = (\xi_{k_1}, \dots, \xi_{k_j})$ is a list of independent column vectors and, $\forall \ell \in R_j := \{k_{j-1} + 1, \dots, k_j - 1\}$, ξ_ℓ can be expressed as a linear combination of the column vectors in K_{j-1} .

By Corollary 6.10, the Gram-Schmidt process determines a unique orthogonal matrix $A_r^* = [u_1^*, u_2^*, \dots, u_r^*] \in \mathbb{R}^{m \times r}$ and a unique upper triangular matrix such that

$$A_r = A_r^* \begin{bmatrix} b_{11} & b_{21} & \dots & b_{r1} \\ & b_{22} & \dots & b_{r2} \\ & & \ddots & \vdots \\ & & & b_{rr} \end{bmatrix}. \quad (6.44)$$

By definition of the column rank of a matrix, we have $r \leq m$.

In the rest of this proof, we insert each column vector in $X = \{\xi_1, \xi_2, \dots, \xi_n\} \setminus \{u_1, u_2, \dots, u_r\}$ back into (6.44) and show that the QR form of (6.44) is maintained. For those zero column vectors in (S-1), we have

$$\begin{aligned} A_\xi &= [\xi_1 \dots \xi_{k_1-1} \ u_1 \ u_2 \ \dots \ u_r] \\ &= A_r^* \begin{bmatrix} 0 & \dots & 0 & b_{11} & b_{21} & \dots & b_{r1} \\ 0 & \dots & 0 & & b_{22} & \dots & b_{r2} \\ \vdots & \ddots & \vdots & & & \ddots & \vdots \\ 0 & \dots & 0 & & & & b_{rr} \end{bmatrix}. \end{aligned} \quad (6.45)$$

For each ξ_ℓ with $\ell \in R_j$ in (S-2), we have

$$\begin{aligned} &[u_1, u_2, \dots, u_{j-1}, \xi_\ell] \\ &= [u_1^*, u_2^*, \dots, u_{j-1}^*] \begin{bmatrix} b_{11} & \dots & b_{j-1,1} & c_{\ell,1} \\ & \ddots & \vdots & \vdots \\ & & b_{j-1,j-1} & c_{\ell,j-1} \end{bmatrix}, \end{aligned} \quad (6.46)$$

where $\xi_\ell = c_{\ell,1}u_1^* + \dots + c_{\ell,j-1}u_{j-1}^*$. With (6.45) as the induction basis and (6.46) as the inductive step, it is straightforward to prove by induction that we have $A = A_r^*R$ where R is an upper triangular matrix.

If $r = m$, Definitions 6.29 and 6.6 complete the proof. Otherwise $r < m$ and the proof is completed by the well-known fact in linear algebra that a list of orthonormal vectors can be extended to an orthonormal basis. \square

Lemma 6.32. An orthogonal matrix preserves the 2-norm of the vectors it acts on.

Proof. Definition 6.29 yields

$$\forall \mathbf{x} \in \text{dom}(Q), \quad \|Q\mathbf{x}\|_2^2 = \mathbf{x}^T Q^T Q \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2. \quad \square$$

Theorem 6.33. Consider an over-determined linear system $A\mathbf{x} = \mathbf{b}$ where $A \in \mathbb{R}^{m \times n}$ and $m \geq n$. The discrete linear least square problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|A\mathbf{x} - \mathbf{b}\|_2^2$$

is solved by \mathbf{x}^* satisfying

$$R_1 \mathbf{x}^* = \mathbf{c}, \quad (6.47)$$

where $R_1 \in \mathbb{R}^{n \times n}$ and $\mathbf{c} \in \mathbb{R}^n$ result from the QR factorization of A :

$$Q^T A = R = \begin{bmatrix} R_1 \\ \mathbf{0} \end{bmatrix}, \quad Q^T \mathbf{b} = \begin{bmatrix} \mathbf{c} \\ \mathbf{r} \end{bmatrix}. \quad (6.48)$$

Furthermore, the minimum is $\|\mathbf{r}\|_2^2$.

Proof. For any $\mathbf{x} \in \mathbb{R}^n$, we have

$$\|A\mathbf{x} - \mathbf{b}\|_2^2 = \|Q^T A\mathbf{x} - Q^T \mathbf{b}\|_2^2 = \|R_1 \mathbf{x} - \mathbf{c}\|_2^2 + \|\mathbf{r}\|_2^2,$$

where the first step follows from Lemma 6.32. \square

6.5 Problems

6.5.1 Theoretical questions

- I. Fill in the details for the proof of Theorem 6.4.
- II. Consider the Chebyshev polynomials of the first kind.
 - (a) Show that they are orthogonal on $[-1, 1]$ with respect to the inner product in Theorem 6.4 with the weight function $\rho(x) = \frac{1}{\sqrt{1-x^2}}$.
 - (b) Normalize the first three Chebyshev polynomials to arrive at an orthonormal system.
- III. Least-square approximation of a continuous function. Approximate the circular arc given by the equation $y(x) = \sqrt{1-x^2}$ for $x \in [-1, 1]$ by a quadratic polynomial with respect to the inner product in Theorem 6.4.
 - (a) $\rho(x) = \frac{1}{\sqrt{1-x^2}}$ with Fourier expansion,
 - (b) $\rho(x) = \frac{1}{\sqrt{1-x^2}}$ with normal equations.
- IV. Discrete least square via orthonormal polynomials. Consider the example on the table of sales record in Example 6.12.
 - (a) Starting from the independent list $(1, x, x^2)$, construct orthonormal polynomials by the Gram-Schmidt process using

$$\langle u(t), v(t) \rangle = \sum_{i=1}^N \rho(t_i) u(t_i) v(t_i) \quad (6.49)$$

as the inner product with $N = 12$ and $\rho(x) = 1$.

- (b) Find the best approximation $\hat{\varphi} = \sum_{i=0}^2 a_i x^i$ such that $\|y - \hat{\varphi}\| \leq \|y - \sum_{i=0}^2 b_i x^i\|$ for all $b_i \in \mathbb{R}$. Verify that $\hat{\varphi}$ is the same as that of the example on the table of sales record in the notes.
- (c) Suppose there are other tables of sales record in the same format as that in the example. Values of N and x_i 's are the same, but the values of y_i 's are different. Which of the above calculations can be reused? Which cannot be reused? What advantage of orthonormal polynomials over normal equations does this reuse imply?

6.5.2 Programming assignments

- A. Write a program to perform discrete least square via normal equations. Your subroutine should take two arrays x and y as the input and output three coefficients a_0, a_1, a_2 that determines a quadratic polynomial as the best fitting polynomial in the sense of least squares with the weight function $\rho = 1$.

Run your subroutine on the following data.

x	0.0	0.5	1.0	1.5	2.0	2.5	3.0
y	2.9	2.7	4.8	5.3	7.1	7.6	7.7
x	3.5	4.0	4.5	5.0	5.5	6.0	6.5
y	7.6	9.4	9.0	9.6	10.0	10.2	9.7
x	7.0	7.5	8.0	8.5	9.0	9.5	10.0
y	8.3	8.4	9.0	8.3	6.6	6.7	4.1

- B. Write a program to solve the previous discrete least square problem via QR factorization. Report the condition number based on the 2-norm of the matrix G in the normal-equation approach and that of the matrix R_1 in the QR-factorization approach, verifying that the former is much larger than the latter.