

Fall 2022

Instructor: Ashwin Pananjady

7750: Mathematical Foundations of Machine Learning

Linear algebra and probability for data analysis

Homework 2

Released: Sep 8

Due: Sep 20, 11:59pm ET

Note. All external sources and collaborators must be acknowledged in your submission. As stated in the syllabus, unauthorized use of previous semester course materials is strictly prohibited in this course.

Objective. To get comfortable with reasoning about Hilbert spaces and linear approximation.

Resources. Lectures, notes, and modules posted on and before Sep 14.

Notation: Capital boldface letters will typically be matrices, and small boldface letter will be vectors.

Problem 1 (Exercises from lecture). 30 points: In this problem, we will walk you through a few proofs of facts that were mentioned in lecture or notes but not explicitly proved. Parts (a-d) are about Cauchy sequences and completeness. Part (e) is about inner products and the so-called parallelogram law. Parts (f-h) help you verify that the space of finite-variance random variables can be viewed as an inner product space.

- (a) Recall that we considered the vector space of real-valued, continuous functions on $[0, 1]$, denoted by $\mathcal{C}[0, 1]$. Equip this space with the standard inner product, with $\langle f, g \rangle := \int_0^1 f(t)g(t)dt$. Also recall the sequence of functions $(f_n)_{n=1}^\infty$ defined as follows

$$f_n(t) = \begin{cases} 0 & \text{if } 0 \leq t \leq \frac{1}{2} - \frac{1}{2n} \\ 2nt + (1 - n) & \text{if } \frac{1}{2} - \frac{1}{2n} < t \leq \frac{1}{2} \\ 1 & \text{if } \frac{1}{2} < t \leq 1. \end{cases}$$

Show that for any finite n , f_n is a continuous function.

Solution: A function $f(x)$ is defined to be continuous at $x = x_0$ if

$$\lim_{x \rightarrow x_0^-} f(x) = \lim_{x \rightarrow x_0^+} f(x) = \lim_{x \rightarrow x_0} f(x)$$

We begin by checking the limit as t approaches $\frac{1}{2} - \frac{1}{2n}$,

$$\lim_{t \rightarrow (\frac{1}{2} - \frac{1}{2n})^-} f_n(t) = 0 \qquad \lim_{t \rightarrow (\frac{1}{2} - \frac{1}{2n})^+} f_n(t) = 2n\left(\frac{1}{2} - \frac{1}{2n}\right) + (1 - n) = 0.$$

Since the limits are equal at $t = (\frac{1}{2} - \frac{1}{2n})^-$ and $t = (\frac{1}{2} - \frac{1}{2n})^+$, $f_n(t)$ is continuous at $t = \frac{1}{2} - \frac{1}{2n}$.

Next, check the limit when t approaches $\frac{1}{2}$ from the left and right,

$$\lim_{t \rightarrow (\frac{1}{2})^-} f_n(t) = 2n \cdot \frac{1}{2} + (1 - n) = 1 \qquad \lim_{t \rightarrow (\frac{1}{2})^+} f_n(t) = 1.$$

The limits are equal and hence, $f_n(t)$ is continuous at $t = \frac{1}{2}$. Hence, $f_n(t)$ is a continuous function on $[0, 1]$.

- (b) Now consider any pair of positive integers $n < m$, and evaluate the norm $\|f_n - f_m\|$. Your expression should be explicit and depend on n and m (recall that the norm is the one induced by the inner product).

Solution: Let $n, m \in \mathbb{N}, n < m$. This means $\frac{1}{2} - \frac{1}{2n} < \frac{1}{2} - \frac{1}{2m}$. The norm is given by,

$$\|f_n - f_m\| := \left(\int_0^1 (f_n - f_m)^2 dt \right)^{\frac{1}{2}}. \quad (1)$$

We first solve (1) by taking the norm squared.

$$\|f_n - f_m\|^2 = \int_0^1 (f_n - f_m)^2 dt.$$

Splitting the integral for the intervals $t \in [0, \frac{1}{2} - \frac{1}{2n}]$, $t \in [\frac{1}{2} - \frac{1}{2n}, \frac{1}{2} - \frac{1}{2m}]$, $t \in [\frac{1}{2} - \frac{1}{2m}, \frac{1}{2}]$, and $t \in [\frac{1}{2}, 1]$ and evaluating gives,

$$\|f_n - f_m\|^2 = \frac{(m - n)^2}{6m^2n}.$$

Thus, $\|f_n - f_m\| = \frac{1}{\sqrt{6n}} \frac{m-n}{m}$.

- (c) Using your expression above, show that the sequence (f_n) is a Cauchy sequence. Conclude that the inner product space introduced above is not complete by thinking about f_∞ .

Solution: To prove f_n is Cauchy, we will show that $\forall \epsilon > 0, \exists N \in \mathbb{N}$, such that $\forall n, m > N$, $\|f_n - f_m\| < \epsilon$.

There are many appropriate choices for selecting such an N , where $n, m > N$. Let us begin with $\epsilon > 0$, and $n, m > N_\epsilon$ where I choose $N = N_\epsilon$ such that $\sqrt{\frac{2}{3N_\epsilon}} < \epsilon$.

$$\begin{aligned} \|f_n - f_m\| &= \frac{1}{\sqrt{6n}} \frac{m-n}{m} = \frac{1}{\sqrt{6n}} - \frac{1}{\sqrt{6n}} \frac{n}{m} < \frac{1}{\sqrt{6n}} + \frac{1}{\sqrt{6n}} \quad (m > n) \\ &= \sqrt{\frac{2}{3N_\epsilon}} < \epsilon. \end{aligned}$$

Thus, f_n is Cauchy.

However, as $n \rightarrow \infty$, f_n becomes a Heaviside step function and is not in $\mathcal{C}[0, 1]$ anymore.

Therefore, the inner product space is not complete.

$$\lim_{n \rightarrow \infty} f_n = \begin{cases} 0, & 0 \leq t < \frac{1}{2} \\ 1, & \frac{1}{2} \leq t \leq 1. \end{cases}$$

- (d) (BONUS:) How might you define the “completion” of the inner product space above so that the resulting space is complete (and hence a Hilbert space)? Note that this will require extra reading that you don’t need to do in principle, just if you’re interested.

- (e) Show that a norm $\|\cdot\|$ defined on a vector space \mathcal{S} is induced by an inner product if and only if the following parallelogram law is satisfied for all $\mathbf{x}, \mathbf{y} \in \mathcal{S}$:

$$2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2 = \|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2.$$

One of these directions is trickier than the other.

Solution:

• Assume $\forall \mathbf{x}, \mathbf{y} \in \mathcal{S}$, the parallelogram law is satisfied. We need to show that the norm $\|\cdot\|$ defined on \mathcal{S} is an induced norm. First express the parallelogram law in terms of an inner product,

$$\begin{aligned}
2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2 &= \|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 \\
\implies 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle + 2\langle \mathbf{x}, \mathbf{y} \rangle) &= \|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 \\
\implies 2\|\mathbf{x} - \mathbf{y}\|^2 + 4\langle \mathbf{x}, \mathbf{y} \rangle &= \|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 \\
\implies \langle \mathbf{x}, \mathbf{y} \rangle &= \frac{1}{4}(\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2).
\end{aligned} \tag{2}$$

It is very clear using (2) that $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$. So symmetry is satisfied.

Next, we show that $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$. Applying the parallelogram law for $\mathbf{x} + \mathbf{z}$ and \mathbf{y} gives

$$\begin{aligned}
2\|\mathbf{x} + \mathbf{z}\|^2 + 2\|\mathbf{y}\|^2 &= \|(\mathbf{x} + \mathbf{z}) + \mathbf{y}\|^2 + \|(\mathbf{x} + \mathbf{z}) - \mathbf{y}\|^2 \\
\implies \|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 &= 2\|\mathbf{x} + \mathbf{z}\|^2 + 2\|\mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y} + \mathbf{z}\|^2.
\end{aligned} \tag{3}$$

Again, applying the parallelogram law for \mathbf{x} and $\mathbf{y} + \mathbf{z}$ gives,

$$\begin{aligned}
2\|\mathbf{x}\|^2 + 2\|\mathbf{y} + \mathbf{z}\|^2 &= \|\mathbf{x} + (\mathbf{y} + \mathbf{z})\|^2 + \|\mathbf{x} - (\mathbf{y} + \mathbf{z})\|^2 \\
\implies \|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 &= 2\|\mathbf{y} + \mathbf{z}\|^2 + 2\|\mathbf{x}\|^2 - \|\mathbf{x} - \mathbf{y} - \mathbf{z}\|^2
\end{aligned} \tag{4}$$

Since the left hand sides are the same for (3) and (4), express,

$$\begin{aligned}
\|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 &= \frac{\underbrace{(\|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2)}_{(3)} + \underbrace{(\|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2)}_{(4)}}{2} \\
&= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 + \|\mathbf{x} + \mathbf{z}\|^2 + \|\mathbf{y} + \mathbf{z}\|^2 - \frac{\|\mathbf{x} - \mathbf{y} + \mathbf{z}\|^2}{2} - \frac{\|\mathbf{x} - \mathbf{y} - \mathbf{z}\|^2}{2}
\end{aligned} \tag{5}$$

Similarly, derive the expression for $\|\mathbf{x} + \mathbf{y} - \mathbf{z}\|^2$, by setting $\mathbf{z} = -\mathbf{z}$, and we get,

$$\|\mathbf{x} + \mathbf{y} - \mathbf{z}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{z}\|^2 + \|\mathbf{y} - \mathbf{z}\|^2 - \frac{\|\mathbf{x} - \mathbf{y} - \mathbf{z}\|^2}{2} - \frac{\|\mathbf{x} - \mathbf{y} + \mathbf{z}\|^2}{2} \tag{6}$$

Subtracting (6) from (5) gives,

$$\begin{aligned}
\|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 - \|\mathbf{x} + \mathbf{y} - \mathbf{z}\|^2 &= \|\mathbf{x} + \mathbf{z}\|^2 - \|\mathbf{x} - \mathbf{z}\|^2 + \|\mathbf{y} + \mathbf{z}\|^2 - \|\mathbf{y} - \mathbf{z}\|^2 \\
\frac{1}{4}\|(\mathbf{x} + \mathbf{y}) + \mathbf{z}\|^2 - \|(\mathbf{x} + \mathbf{y}) - \mathbf{z}\|^2 &= \frac{1}{4}(\|\mathbf{x} + \mathbf{z}\|^2 - \|\mathbf{x} - \mathbf{z}\|^2) + \frac{1}{4}(\|\mathbf{y} + \mathbf{z}\|^2 - \|\mathbf{y} - \mathbf{z}\|^2) \\
\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle &= \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle
\end{aligned}$$

as desired by using (2).

Next, we need to show $\langle \alpha \mathbf{x}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle$, $\forall \alpha \in \mathbb{R}$. Let $\alpha = -1$, it is clear that $\langle -\mathbf{x}, \mathbf{y} \rangle = -\langle \mathbf{x}, \mathbf{y} \rangle$. Since bilinearity holds as shown above, we have $\langle \alpha \mathbf{x}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle$, $\forall \alpha \in \mathbb{Z}$, where \mathbb{Z} is the set of integers. Set $\alpha = \frac{p}{q}$, where $p, q \in \mathbb{Z}$, and $q \neq 0$, and choose $\mathbf{x}' = \frac{\mathbf{x}}{q}$ we have,

$$q \langle \alpha \mathbf{x}, \mathbf{y} \rangle = q \langle \frac{p}{q} \mathbf{x}, \mathbf{y} \rangle = q \langle p \frac{\mathbf{x}}{q}, \mathbf{y} \rangle = q \langle p \mathbf{x}', \mathbf{y} \rangle = p \langle q \mathbf{x}', \mathbf{y} \rangle = p \langle \mathbf{x}, \mathbf{y} \rangle.$$

Dividing by q on both sides gives,

$$\langle \alpha \mathbf{x}, \mathbf{y} \rangle = \frac{p}{q} \langle \mathbf{x}, \mathbf{y} \rangle = \alpha \langle \mathbf{x}, \mathbf{y} \rangle \quad \forall \alpha \in \mathbb{R}$$

We can claim the equality above holds for $\alpha \in \mathbb{R}$ due to the construction of \mathbb{R} by taking a Cauchy sequence of rationals, i.e., for $r \in \mathbb{R}$, $\exists \{q_n\}_n$, $q_i \in \mathbb{Q}$, such that $\lim_{n \rightarrow \infty} q_n = r$.

Finally, observe that $\langle \mathbf{x}, \mathbf{x} \rangle = \frac{1}{4} \|2\mathbf{x}\|^2 = \|\mathbf{x}\|^2$. Therefore, $\sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \|\mathbf{x}\|$. Hence, $\langle \cdot, \cdot \rangle$ is a valid inner product and the norm induced by it is given by $\|\cdot\|^2 = \langle \cdot, \cdot \rangle$.

- Next, assume that $\|\cdot\|^2 = \langle \cdot, \cdot \rangle$, and we wish to prove the parallelogram law. Using $\mathbf{x} + \mathbf{y}$ and $\mathbf{x} - \mathbf{y}$,

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2\langle \mathbf{x}, \mathbf{x} \rangle + 2\langle \mathbf{y}, \mathbf{y} \rangle = 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2.$$

- (f) Consider the space of all real-valued random variables with finite-variance. Show that this is a vector space, and call it \mathcal{S} .

Solution: We need to show that the space of all real-valued random variables with finite-variance is closed under addition and scalar multiplication.

Given $X, Y \in \mathcal{S}$ are real valued random variables with finite variances, $X + Y$ is also a real valued random variable with finite variance, i.e., $\text{Var}(X) + \text{Var}(Y)$ is finite. Because $\mathbb{E}[X] = \int_{\mathcal{X}} x f_X(x)$, so $\mathbb{E}[X] \in \mathbb{R}$. We have $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$. The square of a random variable is also a random variable, so $\mathbb{E}[X^2] \in \mathbb{R}$. Therefore, $\text{Var}[X] \in \mathbb{R}$ is finite, and so is $\text{Var}(X) + \text{Var}(Y)$.

To show closure under scalar multiplication, we need to show $\text{Var}(aX + b) = a^2 \text{Var}(X)$ is finite. Let $a, b \in \mathbb{R}$.

$$\begin{aligned} \text{Var}(aX + b) &= \mathbb{E}[(aX + b - a\mathbb{E}[X] - b)^2] = \mathbb{E}[a^2(X - \mathbb{E}[X])^2] \\ &= a^2 \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= a^2 \text{Var}(X) \end{aligned}$$

Hence, the space of all real-valued random variables with finite variance is a vector space.

- (g) Note that for $X, Y \in \mathcal{S}$, we may view $\mathbb{E}[XY]$ as a function mapping $\mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$. Show that this is a valid inner product on this vector space.

Solution: Given the inner product defined as $\langle X, Y \rangle = \mathbb{E}[XY]$, we need to show symmetricity, bilinearity, and positive definiteness.

It is clear that it is symmetric: $\langle X, Y \rangle = \mathbb{E}[XY] = \mathbb{E}[YX] = \langle Y, X \rangle$.

Next, $\forall a \in \mathbb{R}$, $\langle aX, Y \rangle = \mathbb{E}[aXY] = a\mathbb{E}[XY] = a\langle X, Y \rangle$.

For $X, Y, Z \in \mathcal{S}$, $\langle X + Y, Z \rangle = \mathbb{E}[(X + Y)Z] = \mathbb{E}[XZ + YZ] = \mathbb{E}[XZ] + \mathbb{E}[YZ] = \langle X, Z \rangle + \langle Y, Z \rangle$.

Finally, $\langle X, X \rangle = \mathbb{E}[X^2] = \int_{\mathcal{X}} x^2 f_X(x) dx$. We have $x^2 \geq 0, \forall x \in \mathbb{R}$, and $f_X(x) \geq 0, \forall x \in \mathbb{R}$. Hence, $\mathbb{E}[X^2]$ is always non-negative and is 0 if and only if $X = 0$.

- (h) For any pair of finite variance random variables (X, Y) , the *conditional expectation* $\mathbb{E}[Y|X]$ is a function of X that is known to satisfy the following property: for all functions¹ ϕ

$$\mathbb{E}[(Y - \mathbb{E}[Y|X])\phi(X)] = 0.$$

Using this definition, prove that the mean squared error $\mathbb{E}[(Y - \phi(X))^2]$ of estimating Y from X is minimized by choosing $\phi(X) = \mathbb{E}[Y|X]$. I.e., the conditional expectation minimizes the mean squared error of estimation.

Hint: Think about how we proved the orthogonality principle without necessarily trying to formally define a subspace.

Solution:

$$\begin{aligned} \mathbb{E}[(Y - \phi(X))^2] &= \mathbb{E}[(Y - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - \phi(X))^2] \\ &= \mathbb{E}[\underbrace{(Y - \mathbb{E}[Y|X])^2}_{(a)} + 2\underbrace{(Y - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - \phi(X))}_{(b)} \\ &\quad + \underbrace{(\mathbb{E}[Y|X] - \phi(X))^2}_{(c)}] \end{aligned}$$

Analyzing (a) and (b):

$$(a) = \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] = \mathbb{E}[(Y - \mathbb{E}[Y|X]) \underbrace{(Y - \mathbb{E}[Y|X])}_{\phi_1(X)}] = \mathbb{E}[(Y - \mathbb{E}[Y|X])\phi_1(X)] = 0,$$

by using the property that $\mathbb{E}[(Y - \mathbb{E}[Y|X])\phi_1(X)] = 0$.

Next, we look at (b),

$$(b) = 2\mathbb{E}[(Y - \mathbb{E}[Y|X]) \underbrace{(\mathbb{E}[Y|X] - \phi(X))}_{\phi_2(X)}] = 2\mathbb{E}[(Y - \mathbb{E}[Y|X])\phi_2(X)] = 0.$$

So, we have $\mathbb{E}[(Y - \phi(X))^2] = \mathbb{E}[(\mathbb{E}[Y|X] - \phi(X))^2]$. Therefore, $\mathbb{E}[(Y - \phi(X))^2]$ is minimum when $\mathbb{E}[Y|X] - \phi(X) = 0$, i.e., $\phi(X) = \mathbb{E}[Y|X]$.

¹In reality you need a measurability condition that we will ignore.

Problem 2 (Gram matrices and Gram–Schmidt). 20 points:

- (a) As you know, a square $N \times N$ matrix \mathbf{G} is *invertible* if

$$\mathbf{x}_1 \neq \mathbf{x}_2 \Leftrightarrow \mathbf{G}\mathbf{x}_1 \neq \mathbf{G}\mathbf{x}_2.$$

That is, $\mathbf{G}\mathbf{x}$ is different for every different \mathbf{x} . In other words, if you can show that $\mathbf{G}\mathbf{x} = \mathbf{0}$ only if $\mathbf{x} = \mathbf{0}$, then you have shown that \mathbf{G} is invertible.

Let $\mathbf{v}_1, \dots, \mathbf{v}_N$ be N linearly independent vectors in a Hilbert space, and let $\mathcal{T} = \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_N\}$. Show that if $\mathbf{z} \in \mathcal{T}$ and $\langle \mathbf{v}_n, \mathbf{z} \rangle = 0$ for all $n = 1, \dots, N$, then it must be true that $\mathbf{z} = \mathbf{0}$.

Solution: Since $\mathbf{z} \in \mathcal{T}$, $\mathbf{z} = \sum_{n=1}^N \alpha_n \mathbf{v}_n$ for some α .

$$\langle \mathbf{z}, \mathbf{z} \rangle = \left\langle \sum_{n=1}^N \alpha_n \mathbf{v}_n, \mathbf{z} \right\rangle = \sum_{n=1}^N \alpha_n \langle \mathbf{v}_n, \mathbf{z} \rangle = \sum_{n=1}^N \alpha_n * 0 = 0$$

$\langle \mathbf{z}, \mathbf{z} \rangle = 0$, so $\mathbf{z} = \mathbf{0}$.

- (b) Show that if $\mathbf{v}_1, \dots, \mathbf{v}_N$ are N linearly independent vectors in a Hilbert space, then the Gram matrix

$$\mathbf{G} = \begin{bmatrix} \langle \mathbf{v}_1, \mathbf{v}_1 \rangle & \langle \mathbf{v}_2, \mathbf{v}_1 \rangle & \cdots & \langle \mathbf{v}_N, \mathbf{v}_1 \rangle \\ \langle \mathbf{v}_1, \mathbf{v}_2 \rangle & \langle \mathbf{v}_2, \mathbf{v}_2 \rangle & & \langle \mathbf{v}_N, \mathbf{v}_2 \rangle \\ \vdots & & \ddots & \vdots \\ \langle \mathbf{v}_1, \mathbf{v}_N \rangle & \cdots & & \langle \mathbf{v}_N, \mathbf{v}_N \rangle \end{bmatrix},$$

is invertible. (Hint: use part (a).)

Solution: We will prove with contrapositive, showing that if \mathbf{G} is not invertible, then the $\{\mathbf{v}_n\}$ cannot be linearly independent. Suppose there is an $\mathbf{x} \neq \mathbf{0}$ but $\mathbf{G}\mathbf{x} = \mathbf{0}$. Consider the function

$$\mathbf{z} = \sum_{k=1}^N x_k \mathbf{v}_k,$$

where the coefficients x_k are the entries of \mathbf{x} . Since $\mathbf{G}\mathbf{x} = \mathbf{0}$,

$$\langle \mathbf{z}, \mathbf{v}_n \rangle = \left\langle \sum_{k=1}^N x_k \mathbf{v}_k, \mathbf{v}_n \right\rangle = \sum_{k=1}^N \langle \mathbf{v}_k, \mathbf{v}_n \rangle x_k = (\mathbf{G}\mathbf{x})_n = 0 \quad \text{for all } n = 1, \dots, N.$$

Thus \mathbf{z} is orthogonal to all of the \mathbf{v}_n . But then

$$\langle \mathbf{z}, \mathbf{z} \rangle = \left\langle \sum_{n=1}^N x_n \mathbf{v}_n, \mathbf{z} \right\rangle = \sum_{n=1}^N x_n \langle \mathbf{v}_n, \mathbf{z} \rangle = 0$$

and so, by the definition of inner product, it must be the case that $\mathbf{z} = \mathbf{0}$. This means that

$$\sum_{k=1}^N x_k \mathbf{v}_k = \mathbf{0}$$

for some non-zero vector \mathbf{x} . Thus, $\{\mathbf{v}_n\}$ are not linearly independent.

- (c) Let us now explore an algorithm that takes a basis for a subspace and produces an orthonormal basis for that same subspace. Let $\mathbf{v}_1, \dots, \mathbf{v}_N$ be a basis for a subspace \mathcal{T} of a Hilbert space with inner product $\langle \cdot, \cdot \rangle$ with induced norm $\|\cdot\|$. Define

$$\psi_1 = \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|},$$

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

$$\mathbf{u}_k = \mathbf{v}_k - \sum_{\ell=1}^{k-1} \langle \mathbf{v}_k, \psi_\ell \rangle \psi_\ell,$$

$$\psi_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}.$$

Argue that for the \mathbf{u}_2 produced above that $\|\mathbf{u}_2\| > 0$, and so ψ_2 is well defined.

Solution: By definition $\mathbf{u}_2 = \mathbf{v}_2 - \langle \mathbf{v}_2, \psi_1 \rangle \psi_1$. By the reverse triangle equality we then have

$$\|\mathbf{u}_2\| \geq \|\mathbf{v}_2\| - \|\langle \mathbf{v}_2, \psi_1 \rangle \psi_1\| = \|\mathbf{v}_2\| - |\langle \mathbf{v}_2, \psi_1 \rangle|$$

Since by definition $\|\psi_1\| = 1$. Now since $\|\cdot\|$ is induced by $\langle \cdot, \cdot \rangle$, we can apply the Cauchy-Schwarz inequality:

$$|\langle \mathbf{v}_2, \psi_1 \rangle| \leq \|\mathbf{v}_2\| \|\psi_1\| \Rightarrow \|\mathbf{u}_2\| \geq \|\mathbf{v}_2\| - |\langle \mathbf{v}_2, \psi_1 \rangle| \geq \|\mathbf{v}_2\| - \|\mathbf{v}_2\| \|\psi_1\| = \|\mathbf{v}_2\| - \|\mathbf{v}_2\| = 0$$

Equality of the above only happens when ψ_1 and \mathbf{v}_2 (or equivalently \mathbf{v}_1 and \mathbf{v}_2) are colinear. However, since $\{\mathbf{v}_i\}$ form a basis, they all have to be linearly independent, and thus none of them can be colinear. Thus, equality never happens and $\|\mathbf{u}_2\| > 0$.

- (d) Argue that $\text{span}\{\psi_1, \psi_2\} = \text{span}\{\mathbf{v}_1, \mathbf{v}_2\}$, and show that ψ_1 and ψ_2 are orthonormal. Hence $\{\psi_1, \psi_2\}$ is an orthonormal basis for $\text{span}\{\mathbf{v}_1, \mathbf{v}_2\}$.

Solution: First, it is clear that ψ_1 and ψ_2 are linear combinations of \mathbf{v}_1 and \mathbf{v}_2 . Therefore, any vector that can be spanned by ψ_1 and ψ_2 will also be in the span of \mathbf{v}_1 and \mathbf{v}_2 . Thus, $\text{Span}\{\psi_1, \psi_2\} \subseteq \text{Span}\{\mathbf{v}_1, \mathbf{v}_2\}$.

However, it is just as obvious that \mathbf{v}_1 and \mathbf{v}_2 can be written as linear combinations of $\boldsymbol{\psi}_1$ and $\boldsymbol{\psi}_2$. Thus, $\text{Span}\{\mathbf{v}_1, \mathbf{v}_2\} \subseteq \text{Span}\{\boldsymbol{\psi}_1, \boldsymbol{\psi}_2\}$. Putting these two results together proves $\text{Span}\{\boldsymbol{\psi}_1, \boldsymbol{\psi}_2\} = \text{Span}\{\mathbf{v}_1, \mathbf{v}_2\}$.

Now we need to check that $\boldsymbol{\psi}_1$ and $\boldsymbol{\psi}_2$ are orthonormal. By definition, $\boldsymbol{\psi}_1$ and $\boldsymbol{\psi}_2$ are already normalized, so it follows that $\langle \boldsymbol{\psi}_1, \boldsymbol{\psi}_1 \rangle = 1$ and $\langle \boldsymbol{\psi}_2, \boldsymbol{\psi}_2 \rangle = 1$. To check orthogonality, consider that $\boldsymbol{\psi}_2$ is a scaled version of \mathbf{u}_2 , so if $\boldsymbol{\psi}_1$ is orthogonal to \mathbf{u}_2 , it is orthogonal to $\boldsymbol{\psi}_2$.

$$\langle \mathbf{u}_2, \boldsymbol{\psi}_1 \rangle = \langle \mathbf{v}_2 - \langle \mathbf{v}_2, \boldsymbol{\psi}_1 \rangle \boldsymbol{\psi}_1, \boldsymbol{\psi}_1 \rangle = \langle \mathbf{v}_2, \boldsymbol{\psi}_1 \rangle - \langle \mathbf{v}_2, \boldsymbol{\psi}_1 \rangle \langle \boldsymbol{\psi}_1, \boldsymbol{\psi}_1 \rangle = \langle \mathbf{v}_2, \boldsymbol{\psi}_1 \rangle - \langle \mathbf{v}_2, \boldsymbol{\psi}_1 \rangle * 1 = 0$$

Thus, $\boldsymbol{\psi}_1$ and $\boldsymbol{\psi}_2$ are orthonormal.

- (e) Use induction to show that $\{\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_N\}$ is an orthonormal basis for \mathcal{T} . Part of this argument will be to ensure that $\mathbf{u}_k \neq \mathbf{0}$.

Solution: In part (d) we proved the base case. Now we do the induction step. Let us assume $\{\boldsymbol{\psi}_k\}_{k=1}^m$ are orthogonal, and we want to show $\{\boldsymbol{\psi}_k\}_{k=1}^{m+1}$ are also orthogonal.

$$\begin{aligned} \langle \mathbf{u}_{m+1}, \boldsymbol{\psi}_j \rangle &= \langle \mathbf{v}_{m+1} - \sum_{\ell=1}^m \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_\ell \rangle \boldsymbol{\psi}_\ell, \boldsymbol{\psi}_j \rangle \quad \forall j = 1, 2, \dots, m \\ &= \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_j \rangle - \sum_{\ell=1}^m \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_\ell \rangle \langle \boldsymbol{\psi}_\ell, \boldsymbol{\psi}_j \rangle \quad \forall j = 1, 2, \dots, m \\ &= \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_j \rangle - \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_j \rangle * 1 - \sum_{\ell \neq j} \langle \mathbf{v}_{m+1}, \boldsymbol{\psi}_\ell \rangle * 0 = 0 \quad \forall j = 1, 2, \dots, m \end{aligned}$$

Thus $\{\boldsymbol{\psi}_k\}_{k=1}^{m+1}$ are orthogonal, and since by definition they are normalized, they are orthonormal.

Finally, the same logic from part (d) applies here too. We know that $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_N$ are linear combinations of $\mathbf{v}_1, \dots, \mathbf{v}_N$, so $\text{Span}(\{\boldsymbol{\psi}_k\}_{k=1}^N) \subseteq \mathcal{T}$. Put together with the (just proven) fact that $\{\boldsymbol{\psi}_k\}_{k=1}^N$ are orthogonal, this implies that $\text{Span}(\{\boldsymbol{\psi}_k\}_{k=1}^N) = \mathcal{T}$. Thus, $\{\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_N\}$ is an orthonormal basis for \mathcal{T} .

Problem 3. (Linear approximation with “bump” functions). 30 points: In this problem, we will develop the computational framework for approximating a function on $[0, 1]$ using scaled and shifted version of the classic bell-curve bump:

$$\phi(t) = e^{-t^2}.$$

Fix an integer $N > 0$ and define $\phi_k(t)$ as

$$\phi_k(t) = \phi\left(\frac{t - (k - 1/2)/N}{1/N}\right) = \phi(Nt - k + 1/2)$$

for $k = 1, 2, \dots, N$. The $\{\phi_k(t)\}$ are a basis for the subspace

$$\mathcal{T}_N = \text{span}\{\phi_k\}_{k=1}^N.$$

- (a) For a fixed value of N , we can plot all of the $\phi_k(t)$ on the same set of axes. Do this for $N = 10$ and $N = 25$ and one more value of N (of your choosing) and turn in your plots.

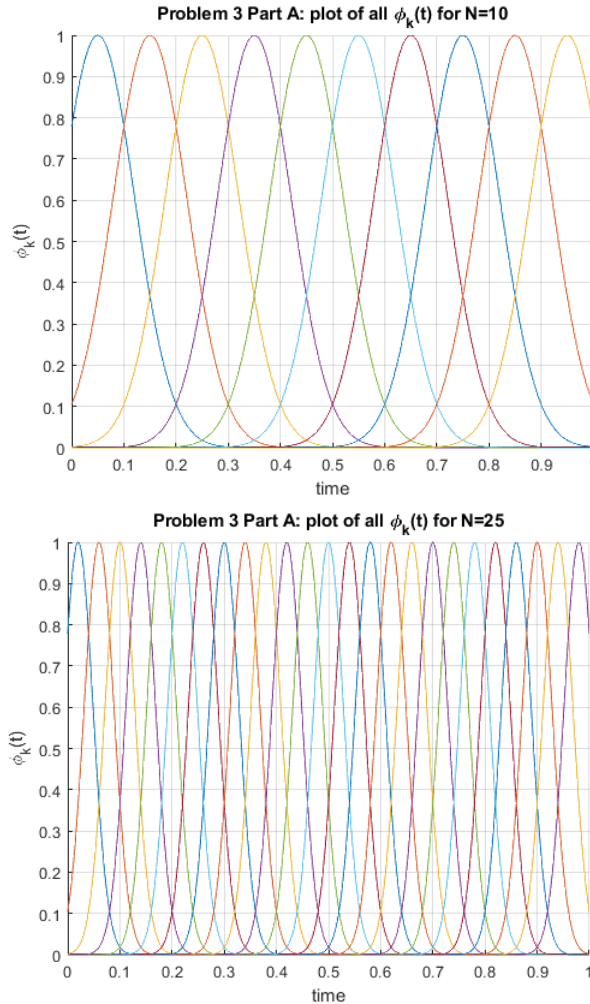


Figure 1: Problem 3, Part A: $\phi_k(t)$ for $N = 10$ and $N = 25$

(b) Since $\{\phi_k\}$ is a basis for \mathcal{T}_N , we can write any $\mathbf{y} \in \mathcal{T}_N$ as

$$y(t) = \sum_{k=1}^N a_k \phi_k(t)$$

for some set of coefficients $a_1, \dots, a_N \in \mathbb{R}^N$. Do this for $N = 4$, and $a_1 = -1/2, a_2 = 3, a_3 = 2, a_4 = -1$ and submit a plot.

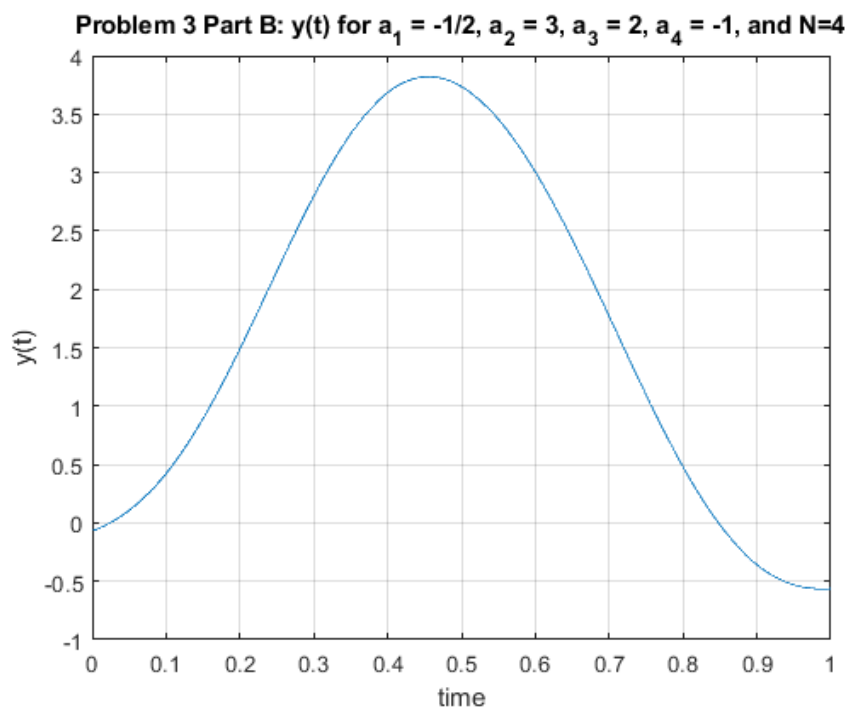


Figure 2: Problem 3, Part B: $y(t)$ for $N = 4$, $\mathbf{a} = [-1/2, 3, 2, -1]^T$

(c) Define the function $f(t)$ on $[0, 1]$ as

$$f(t) = \begin{cases} 4t & 0 \leq t < 1/4 \\ -4t + 2 & 1/4 \leq t < 1/2 \\ -\sin(14\pi t) & 1/2 \leq t \leq 1 \end{cases}$$

Write a function that finds the closest point $\hat{\mathbf{f}}$ in \mathcal{T}_N to \mathbf{f} for any fixed N . By “closest point”, we mean that $\hat{x}(t)$ is the solution to

$$\text{minimize}_{\mathbf{y} \in \mathcal{T}_N} \|\mathbf{f} - \mathbf{y}\|_{L_2([0,1])}, \quad \|\mathbf{f} - \mathbf{y}\|_{L_2([0,1])}^2 = \int_0^1 |f(t) - y(t)|^2 dt.$$

Turn in your code and four plots; one of which has $f(t)$ and $\hat{f}(t)$ plotted on the same set of axes for $N = 5$, and then repeat for $N = 10, 20$, and 50 .

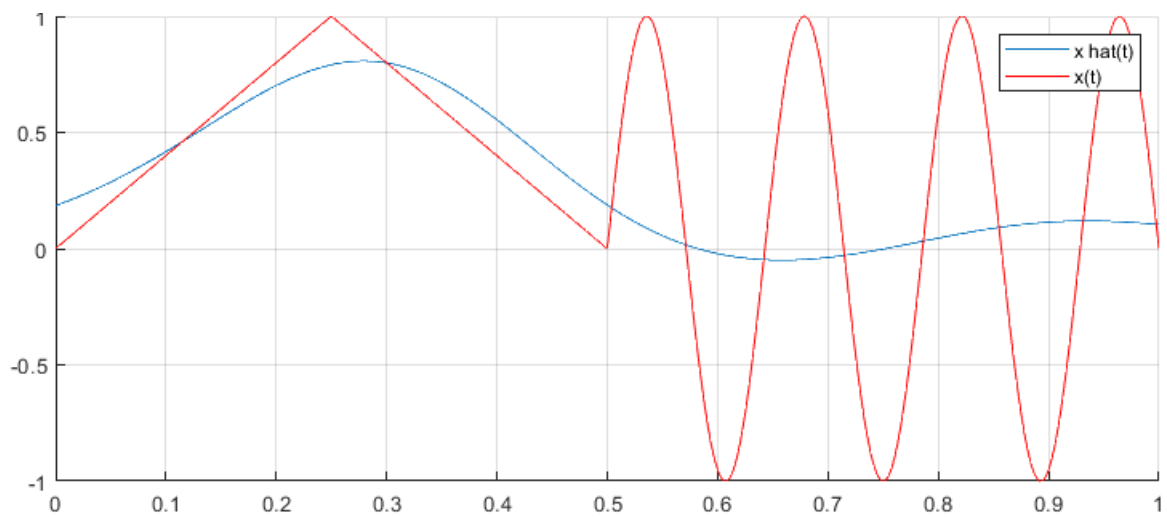


Figure 3: Problem 3, Part C: $\hat{x}(t)$ and $x(t)$ for $N = 5$

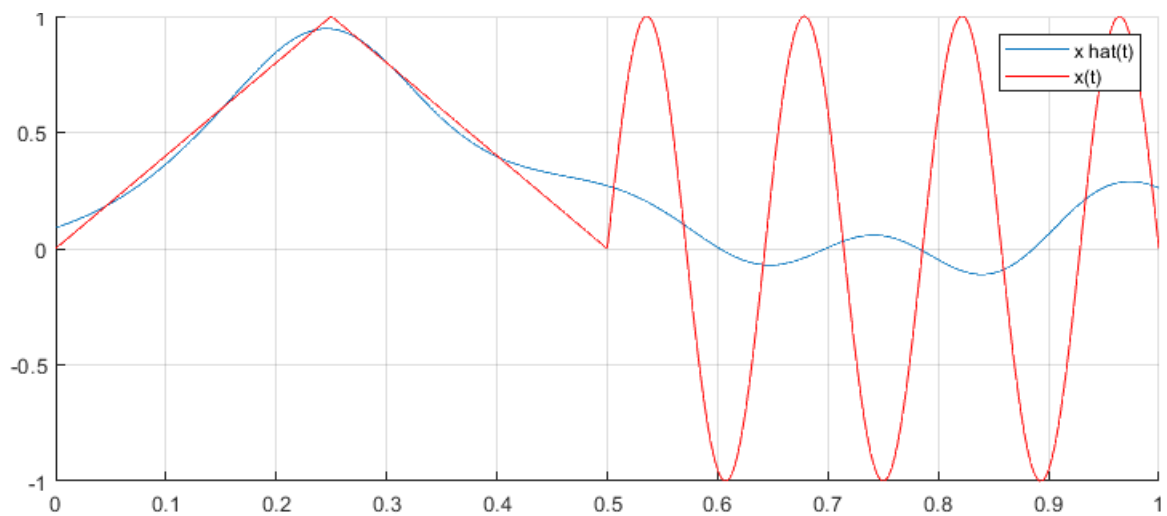


Figure 4: Problem 3, Part C: $\hat{x}(t)$ and $x(t)$ for $N = 10$

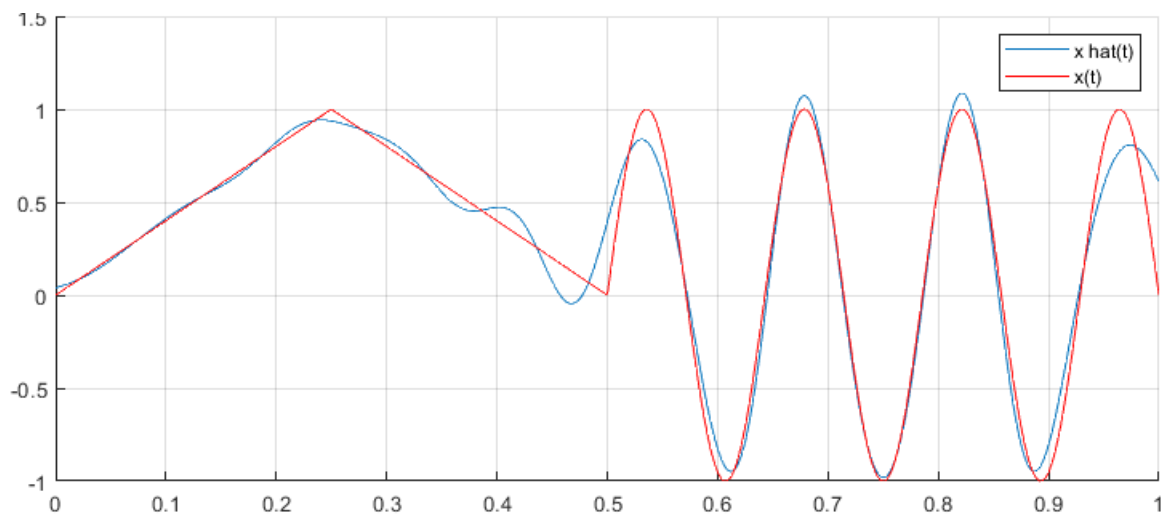


Figure 5: Problem 3, Part C: $\hat{x}(t)$ and $x(t)$ for $N = 20$

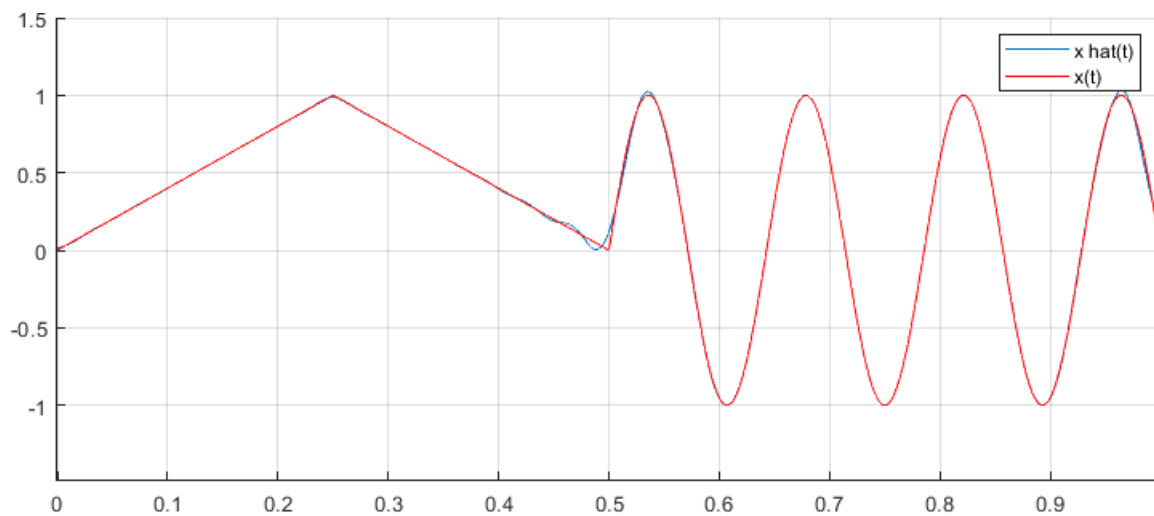


Figure 6: Problem 3, Part C: $\hat{x}(t)$ and $x(t)$ for $N = 50$

```

phi = @(z) exp(-z.^2);
t = linspace(0, 1, 1000);

x = @(z) (z < 1/4).*(4*z) + (z>=1/4).*(z<1/2).*(-4*z+2) ...
- (z>=1/2).*sin(14*pi*z);

N = 50; % N = 25;
G = zeros(N);
b = zeros(N,1);
for jj = 1:N
for kk = 1:N
G(jj,kk) = integral(@(z) phi(N*z - jj + 1/2).*phi(N*z - kk + 1/2),0,1);
end
b(jj) = integral(@(z) phi(N*z - jj + 1/2).*x(z), 0, 1);
end
a = G\b;

xhat = zeros(size(t));
for jj = 1:N
xhat = xhat + a(jj)*phi(N*t - jj + 1/2);
end
figure
hold on
plot(t, xhat);
plot(t, x(t), 'r')
title(sprintf('Problem 3 Part C: plot of x hat(t) for N = %d',N));
legend('x hat(t)', 'x(t)');

```

Problem 4. (Finite dimensional linear regression to predict disease progression). 20 points In this problem, you will run linear regressions on a data set (`diabetes.csv`) containing $d = 10$ predictors (including 7 blood serum measurements) of a response `prog` (diabetes progression) in $n = 442$ patients. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ denote the matrix formed by stacking all these predictors, and $\mathbf{y} \in \mathbb{R}^n$ denote the corresponding responses. You may use any Python package for this problem, but just doing standard linear algebra operations with numpy should suffice.

- (a) Import the data and *standardize* each predictor, i.e., when you look at a particular predictor on all the samples and view it as a vector (i.e. a column vector of \mathbf{X}), you want its empirical mean to be 0 and its empirical variance to be 1.

```
%% Load dataset
T = readtable('diabetes.csv');
N = length(T.prog);
d = 10;

%% Training input and label
X = table2array(T(:,1:d));    Y = T.prog;

%% (a) Standardize input data
X = (X - mean(X))./std(X);
```

- (b) Perform linear regression on this data by writing down the least squares solution **with the intercept**.

```
%% (b) Linear Regression with the intercept
Xi = [ones(N,1) X]; % dataset with column of 1s
what_1 = (Xi'*Xi)\Xi'*Y;
```

The coefficients are,

$$\hat{\mathbf{w}}_1 = \begin{bmatrix} 152.13 \\ -0.46265 \\ -11.386 \\ 24.722 \\ 15.435 \\ -37.476 \\ 22.479 \\ 4.7576 \\ 8.4402 \\ 35.686 \\ 3.236 \end{bmatrix}.$$

The value 152.13 corresponds to the intercept.

- (c) Perform linear regression on the data **without an intercept**. Compare your solution to the previous part and explain what just happened. We would like a justification rooted in linear algebra arguments.

```
%% (c) Linear Regression without the intercept
what_2 = (X'*X)\X'*Y;
```

The coefficients are,

$$\hat{\mathbf{w}}_2 = \begin{bmatrix} -0.46265 \\ -11.386 \\ 24.722 \\ 15.435 \\ -37.476 \\ 22.479 \\ 4.7576 \\ 8.4402 \\ 35.686 \\ 3.236 \end{bmatrix}$$

Only the intercept is missing, but the rest of the coefficients remain the same.

In linear regression, the residual error is given by: $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$, where $\hat{\mathbf{y}}$ is the projection of \mathbf{y} on the linear space spanned by the columns of \mathbf{X} . And we know that \mathbf{e} is orthogonal to any of the columns in \mathbf{X} , since it is the projection of \mathbf{y} onto the span of \mathbf{X} . By introducing a column of $\mathbf{1}$ in \mathbf{X} , and let's call it \mathbf{X}_1 , we make sure that the residual error is also orthogonal to this column of ones in \mathbf{X}_1 . That means, the dot product of \mathbf{e} and $\mathbf{1}$ must be zero, that means the mean of \mathbf{e} must be zero.

- (d) Let us now explore how to obtain a purely data-dependent estimate of the *test error*. Split your data set (at random) into two equal portions, use the first half to fit a linear regression as in part (b), and use the second half to measure the *test* error in predicting their corresponding response values. Do this 100 times and compute the average test error.

```
%% (d) Split the dataset and compute MSE/RMSE with intercept
idx_train = randperm(N, N/2);
idx_test = setxor(1:N, idx_train);

% One-half as training set
X_train = X(idx_train,:);      Y_train = Y(idx_train);

% One-half as testing set
X_test = X(idx_test,:);        Y_test = Y(idx_test);
```



```
% Linear Reg. with intercept on training set using library function
mdl = fitlm(X_train, Y_train);
trials = 100;
err_arr = zeros(trials,1);
for i = 1:trials
    Y_pred = predict(mdl, X_test);
    err_arr(i) = immse(Y_test, Y_pred);
end
MSE = (mean(err_arr));
```

The root MSE is 58.923.

- (e) Repeat part (d) with the linear regressions computed without intercepts.

```
%% (e) Compute RMSE without intercept
% Linear Reg. without intercept on training set
mdl = fitlm(X_train, Y_train, 'intercept', false);
err_arr = zeros(trials,1);
for i = 1:trials
    Y_pred = predict(mdl, X_test);
    err_arr(i) = immse(Y_test, Y_pred);
end
MSE = (mean(err_arr));
```

The root MSE is 163.59.

- (f) Depending on your observations in the previous two parts, which linear model—with or without intercepts—do you think provides a better representation of this data set?

By comparing the MSE or root MSE in parts (d) and (e), we see that the linear model with the intercept provides a better fit than the model without any intercept.