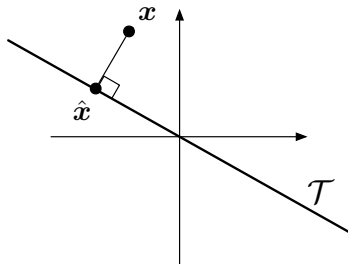


Linear approximation in a Hilbert space

Consider the following problem:

Let \mathcal{S} be a Hilbert space, and let \mathcal{T} be a subspace of \mathcal{S} . Given a $\mathbf{x} \in \mathcal{S}$, what is the **closest point** $\hat{\mathbf{x}} \in \mathcal{T}$?



In other words, find $\hat{\mathbf{x}} \in \mathcal{T}$ that minimizes $\|\mathbf{x} - \hat{\mathbf{x}}\|$; given \mathbf{x} , we want to solve the following optimization program

$$\underset{\mathbf{y} \in \mathcal{T}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{y}\|, \quad (1)$$

where the norm above is the one induced by the inner product. This problem has a unique solution which is characterized by the **orthogonality principle**.

Theorem: Let \mathcal{S} be a Hilbert space, and let \mathcal{T} be a finite dimensional subspace¹. Given an arbitrary $\mathbf{x} \in \mathcal{S}$,

1. there is exactly one $\hat{\mathbf{x}} \in \mathcal{T}$ such that

$$\mathbf{x} - \hat{\mathbf{x}} \perp \mathcal{T}, \quad (2)$$

meaning $\langle \mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} \rangle = 0$ for all $\mathbf{y} \in \mathcal{T}$, and

¹The same results hold when \mathcal{T} is infinite dimensional and is *closed*. We do not prove the infinite dimensional case just because it requires some analysis of infinite sequences which, while not really that difficult, kind of distract from the overall geometrical picture we are trying to paint here.

2. this $\hat{\mathbf{x}}$ is the closest point in \mathcal{T} to \mathbf{x} ; that is, $\hat{\mathbf{x}}$ is the unique minimizer to (1).

Proof: We will show that the first part is true in the next section of the notes, where we show how to explicitly calculate such an $\hat{\mathbf{x}}$.

For the second part, let $\hat{\mathbf{x}}$ be the vector which obeys

$$\hat{\mathbf{e}} = \mathbf{x} - \hat{\mathbf{x}} \perp \mathcal{T}.$$

Let \mathbf{y} be any other vector in \mathcal{T} , and set

$$\mathbf{e} = \mathbf{x} - \mathbf{y}.$$

We will show that

$$\|\mathbf{e}\| > \|\hat{\mathbf{e}}\| \quad (\text{i.e. that } \|\mathbf{x} - \mathbf{y}\| > \|\mathbf{x} - \hat{\mathbf{x}}\|).$$

Note that

$$\begin{aligned} \|\mathbf{e}\|^2 &= \|\mathbf{x} - \mathbf{y}\|^2 = \|\hat{\mathbf{e}} - (\mathbf{y} - \hat{\mathbf{x}})\|^2 \\ &= \langle \hat{\mathbf{e}} - (\mathbf{y} - \hat{\mathbf{x}}), \hat{\mathbf{e}} - (\mathbf{y} - \hat{\mathbf{x}}) \rangle \\ &= \|\hat{\mathbf{e}}\|^2 + \|\mathbf{y} - \hat{\mathbf{x}}\|^2 - \langle \hat{\mathbf{e}}, \mathbf{y} - \hat{\mathbf{x}} \rangle - \langle \mathbf{y} - \hat{\mathbf{x}}, \hat{\mathbf{e}} \rangle. \end{aligned}$$

Since $\mathbf{y} - \hat{\mathbf{x}} \in \mathcal{T}$ and $\hat{\mathbf{e}} \perp \mathcal{T}$,

$$\langle \hat{\mathbf{e}}, \mathbf{y} - \hat{\mathbf{x}} \rangle = 0, \quad \text{and} \quad \langle \mathbf{y} - \hat{\mathbf{x}}, \hat{\mathbf{e}} \rangle = 0,$$

and so

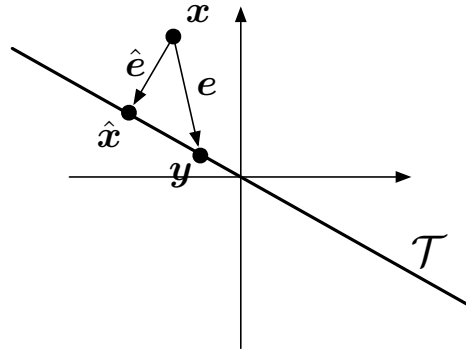
$$\|\mathbf{e}\|^2 = \|\hat{\mathbf{e}}\|^2 + \|\mathbf{y} - \hat{\mathbf{x}}\|^2.$$

Since all three quantities in the expression above are positive and

$$\|\mathbf{y} - \hat{\mathbf{x}}\| > 0 \quad \Leftrightarrow \quad \mathbf{y} \neq \hat{\mathbf{x}},$$

we see that

$$\mathbf{y} \neq \hat{\mathbf{x}} \quad \Leftrightarrow \quad \|\mathbf{e}\| > \|\hat{\mathbf{e}}\|.$$



Computing the best approximation

The orthogonality principle gives us a concrete procedure for actually **computing** the optimal point $\hat{\mathbf{x}}$.

Let N be the dimension of \mathcal{T} , and let $\mathbf{v}_1, \dots, \mathbf{v}_N$ be a basis for \mathcal{T} . We want to find coefficients $a_1, \dots, a_N \in \mathbb{C}$ such that

$$\hat{\mathbf{x}} = a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_N \mathbf{v}_N.$$

The orthogonality principle tells us that

$$\langle \mathbf{x} - \hat{\mathbf{x}}, \mathbf{v}_n \rangle = 0 \quad \text{for } n = 1, \dots, N.$$

This means the a_n must obey

$$\langle \mathbf{x} - \sum_{k=1}^N a_k \mathbf{v}_k, \mathbf{v}_n \rangle = 0 \quad \text{for } n = 1, \dots, N,$$

or moving things around,

$$\langle \mathbf{x}, \mathbf{v}_n \rangle = \sum_{k=1}^N a_k \langle \mathbf{v}_k, \mathbf{v}_n \rangle \quad \text{for } n = 1, \dots, N.$$

Since \mathbf{x} and the $\{\mathbf{v}_n\}$ are given, we know both the $\langle \mathbf{x}, \mathbf{v}_n \rangle$ and the $\langle \mathbf{v}_k, \mathbf{v}_n \rangle$. We are left with a set of N **linear equations** with N unknowns:

$$\begin{bmatrix} \langle \mathbf{v}_1, \mathbf{v}_1 \rangle & \langle \mathbf{v}_2, \mathbf{v}_1 \rangle & \cdots & \langle \mathbf{v}_N, \mathbf{v}_1 \rangle \\ \langle \mathbf{v}_1, \mathbf{v}_2 \rangle & \langle \mathbf{v}_2, \mathbf{v}_2 \rangle & & \langle \mathbf{v}_N, \mathbf{v}_2 \rangle \\ \vdots & & \ddots & \vdots \\ \langle \mathbf{v}_1, \mathbf{v}_N \rangle & \cdots & & \langle \mathbf{v}_N, \mathbf{v}_N \rangle \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{v}_1 \rangle \\ \langle \mathbf{x}, \mathbf{v}_2 \rangle \\ \vdots \\ \langle \mathbf{x}, \mathbf{v}_N \rangle \end{bmatrix}$$

The matrix on the left hand side above is called the **Gram matrix** or **Grammian** of the basis $\{\mathbf{v}_n\}$.

In more compact notation, we want to find $\mathbf{a} \in \mathbb{C}^N$ such that

$$\mathbf{G}\mathbf{a} = \mathbf{b},$$

where $b_n = \langle \mathbf{x}, \mathbf{v}_n \rangle$ and $G_{k,n} = \langle \mathbf{v}_n, \mathbf{v}_k \rangle$.

Two notes on the structure of \mathbf{G} :

- \mathbf{G} is guaranteed to be invertible because the $\{\mathbf{v}_n\}$ are linearly independent. We can comfortably write

$$\mathbf{a} = \mathbf{G}^{-1}\mathbf{b}.$$

- \mathbf{G} is **conjugate symmetric** (“Hermitian”):

$$\mathbf{G} = \mathbf{G}^{\text{H}},$$

where \mathbf{G}^{H} is the conjugate transpose of \mathbf{G} (take the transpose, then take the complex conjugate of all the entries). This fact has algorithmic implications when it comes time to actually solve the system of equations.

Uniqueness

It should be clear that if $\langle \mathbf{e}, \mathbf{v}_k \rangle = 0$ for all of the basis vectors $\mathbf{v}_1, \dots, \mathbf{v}_N$, then $\langle \mathbf{e}, \mathbf{y} \rangle = 0$ for all $\mathbf{y} \in \mathcal{T}$. The converse is also true: if $\langle \mathbf{e}, \mathbf{y} \rangle \neq 0$ for some $\mathbf{y} \in \mathcal{T}$ not equal to $\mathbf{0}$, then $\langle \mathbf{e}, \mathbf{v}_k \rangle \neq 0$ for at least one of the \mathbf{v}_k .

With the work above, this means that a necessary and sufficient condition for $\langle \mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} \rangle = 0$ for all $\mathbf{y} \in \mathcal{T}$ is to have

$$\hat{\mathbf{x}} = \sum_{n=1}^N a_n \mathbf{v}_n, \quad \text{where } \mathbf{a} \text{ satisfies } \mathbf{G}\mathbf{a} = \mathbf{b}.$$

Since \mathbf{G} is square and invertible, there is exactly one such \mathbf{a} , and hence exactly one $\hat{\mathbf{x}}$ that obeys the condition

$$\mathbf{x} - \hat{\mathbf{x}} \perp \mathcal{T}.$$

Example: Let

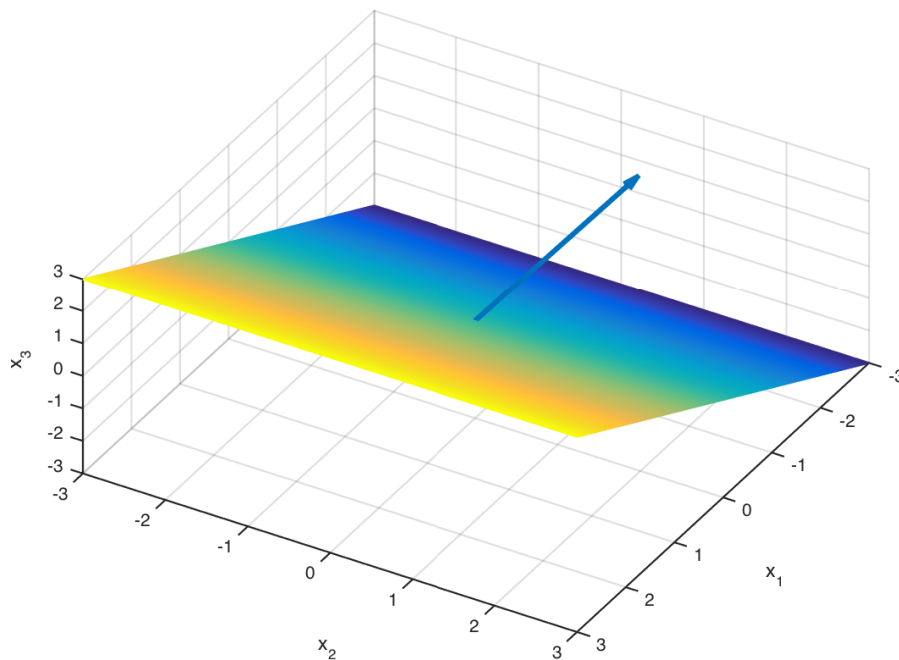
$$\mathcal{T} = \text{Span} \left(\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \right), \quad \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 3 \end{bmatrix}$$

Find the solution to

$$\underset{\mathbf{y} \in \mathcal{T}}{\text{minimize}} \|\mathbf{x} - \mathbf{y}\|_2.$$

(Recall that $\|\cdot\|_2$ in \mathbb{R}^3 is induced by the standard inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{n=1}^3 x_n y_n$.)

Here is a plot of \mathcal{T} and \mathbf{x} :



Solution: We have

$$\mathbf{G} = \qquad \qquad \mathbf{b} =$$

and so

$$\mathbf{G}^{-1} =$$

This means that

$$\mathbf{a} =$$

from which we synthesize the answer

$$\hat{\mathbf{x}} =$$

We can also check that “the error is orthogonal to the approximation”

$$\langle \mathbf{x} - \hat{\mathbf{x}}, \hat{\mathbf{x}} \rangle =$$

Example: Polynomial approximation of e^t

1. We want to calculate a quadratic approximation of $x(t) = e^t$ over the interval $[0, 1]$.
 - (a) Sketch $x(t)$ on $[0, 1]$.
 - (b) One way to do the approximation is to truncate the Taylor expansion

$$e^t = 1 + t + \frac{t^2}{2} + \frac{t^3}{6} + \frac{t^4}{24} + \cdots$$

In this case the approximation is

$$\tilde{x}_{\text{taylor}}(t) = 1 + t + t^2/2.$$

Sketch $\tilde{x}_{\text{taylor}}(t)$.

- (c) Is the truncated Taylor approximation the best possible quadratic approximation? The answer is no if we are interested in minimizing the energy of the error. We want to find the second order polynomial

$$\tilde{x}(t) = a_1 + a_2 t + a_3 t^2$$

that minimizes

$$\|e^t - \tilde{x}\|_{L_2([0,1])} = \sqrt{\int_0^1 |e^t - \tilde{x}(t)|^2 dt}.$$

We set this up as a subspace approximation problem. Set

$$v_1(t) = 1, \quad v_2(t) = t, \quad v_3(t) = t^2,$$

and set $\mathcal{T} = \text{Span}(\{v_1, v_2, v_3\})$. We of course now know a systematic way to find the best $\tilde{x} \in \mathcal{T}$. Start by calculating the Gram matrix \mathbf{G} (recall that $G_{ij} = \langle v_i, v_j \rangle$).

- (d) Calculate the right-hand-side \mathbf{b} (recall that $b_i = \langle x, v_i \rangle$). Write down the system of equations that need to be solved. You can use MATLAB to solve this system.
 - (e) Calculate the error for both the Taylor approximation and the optimal approximation computed above. Plot $x(t)$, $\tilde{x}_{\text{taylor}}(t)$, and $\tilde{x}(t)$ on the same set of axes.

2. We consider the same problem, but calculate the error in a different way. The norm we use to measure the error is

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_S^2 = \int_0^1 w(t)|x(t) - \tilde{x}(t)|^2 dt$$

where

$$w(t) = 16(t - 1/2)^2.$$

- (a) Plot $w(t)$ and argue that measuring the error in this way will penalize mismatch at the ends of the interval than in the middle.
- (b) Write down the inner product that induces $\|\cdot\|_S$.
- (c) Find the second order polynomial that is the best approximation to e^t in the $\|\cdot\|_S$ norm.