

Prelude to orthobases: moving to infinite dimensions

The last set of notes gave us a computationally concrete technique for *approximating* an arbitrary function \mathbf{f} from a Hilbert space using a linear combination of *fixed* elements \mathbf{v}_n from the space. The best approximation (in terms of the norm induced by the inner product in the Hilbert space) to \mathbf{f} is given by

$$\hat{\mathbf{f}}(t) = \sum_{n=1}^N \hat{\alpha}_n v_n(t), \quad \hat{\alpha} = \mathbf{G}^{-1} \mathbf{b}, \quad (1)$$

where

$$G_{i,j} = \langle \mathbf{v}_j, \mathbf{v}_i \rangle, \quad b_i = \langle \mathbf{f}, \mathbf{v}_i \rangle.$$

This is a well-posed problem in finite dimensions since we know that \mathbf{G} is invertible when the $\{\mathbf{v}_n\}$ are linearly independent. If $\mathbf{f} \in \text{Span}(\{\mathbf{v}_1, \dots, \mathbf{v}_N\})$, then $\hat{\mathbf{f}} = \mathbf{f}$, and (1) can be viewed as a “reproducing formula” for \mathbf{f} .

What we would like to do is take this idea to infinite dimensions. That is, we would like an infinite sequence of vectors $\mathbf{v}_1, \mathbf{v}_2, \dots$ so that for any fixed $\mathbf{f} \in \mathcal{H}$, when we create the approximation

$$\hat{\mathbf{f}}_N = \text{best approximation to } \mathbf{f} \text{ in } \text{Span}(\{\mathbf{v}_1, \dots, \mathbf{v}_N\})$$

using (1) above, $\hat{\mathbf{f}}_N$ gets closer and closer to \mathbf{f} as N increases. That is, we would like

$$\lim_{N \rightarrow \infty} \|\hat{\mathbf{f}}_N - \mathbf{f}_N\| = 0,$$

or equivalently $\lim_{N \rightarrow \infty} \hat{\mathbf{f}}_N = \mathbf{f}$. As N increases, the space spanned by the $\{\mathbf{v}_n\}_{n=1}^N$ is becoming richer and richer, eventually encompass-

ing the entire Hilbert space. In the limit, we would like to write

$$f(t) = \sum_{n=1}^{\infty} \alpha_n \mathbf{v}_n(t),$$

i.e. represent \mathbf{f} by the (infinite) list of numbers $\{\alpha_n\}$.

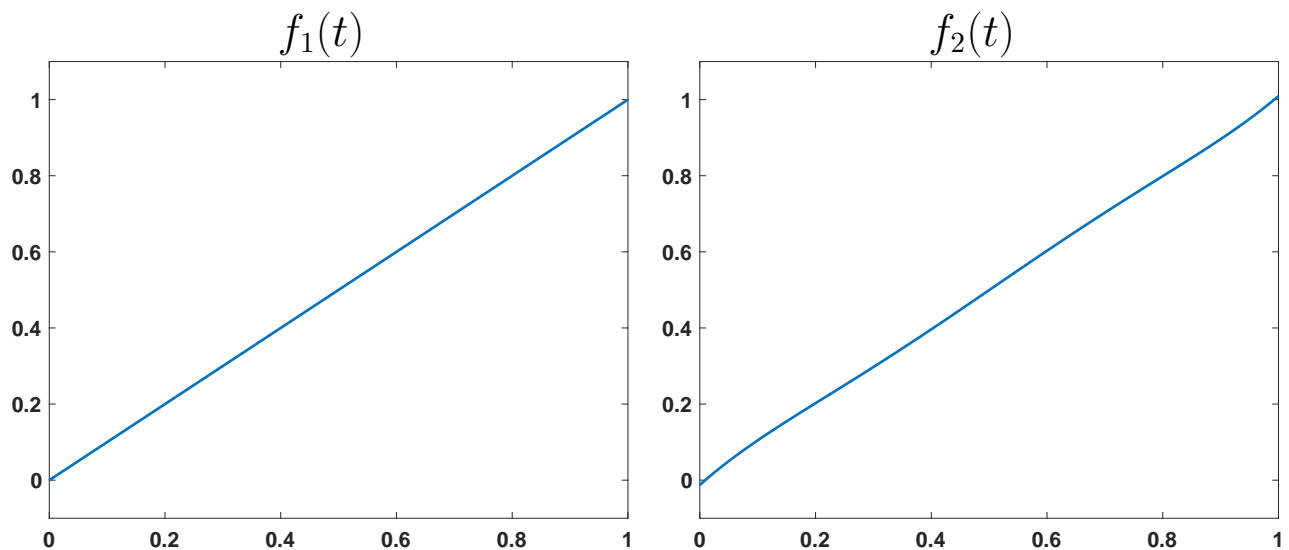
This process, however, is not always as straightforward as it sounds. Things can go wrong, even in familiar, well-defined cases.

The example below illustrates the problem clearly in one such familiar setting. On the left, we plot

$$f_1(t) = t,$$

while on the right, we plot

$$f_2(t) = 2.716 t^5 - 6.898 t^4 + 6.255 t^3 - 2.4068 t^2 + 1.356 t - 0.013$$



Although they have very different expansion coefficients, these functions are almost exactly the same. In fact, the relative error in the

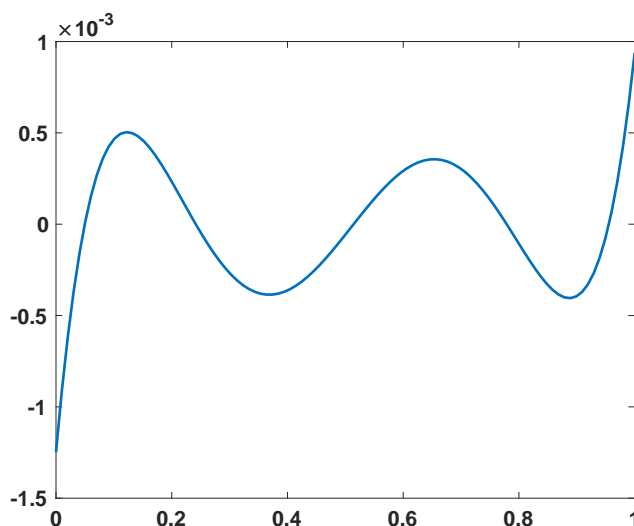
L_2 norm is

$$\frac{\|\mathbf{f}_1 - \mathbf{f}_2\|_2}{\|\mathbf{f}_1\|} \approx 3.25 \cdot 10^{-5}.$$

One conclusion we can draw from this is while the vectors $\{1, t, t^2, t^3, t^4, t^5\}$ are technically linearly independent, functionally they are not — we can build up very close approximations of members of that set through linear combinations of the other elements. Equivalently, there is a good 5th-order approximation of the $\mathbf{0}$ function on $[0, 1]$; we have

$$f_2(t) - f_1(t) = 2.716 t^5 - 6.898 t^4 + 6.255 t^3 - 2.4068 t^2 + 0.356 t - 0.013,$$

which is plotted below,



Notice that the scale of this plot is 10^{-3} , i.e. three orders of magnitude below the plots above (which is why f_1 looks pretty similar to f_2 in those plots).

For polynomials, this effect becomes far more dramatic as we increase their order. While we were able to get within $1.5 \cdot 10^{-3}$ pointwise with a 5th order polynomial, we can get within $1.5 \cdot 10^{-6}$ with a 9th order polynomial, and within $6 \cdot 10^{-8}$ with a 11th order polynomial.

So while we might be able to say things like $f(t) \approx \sum_{n=0}^N \alpha_n t^n$ for large enough N , we do not really have confidence in the $\{\alpha_n\}$ as a representation for \mathbf{f} , as the mapping from the function \mathbf{f} to the list of numbers $\{\alpha_n\}$ is unstable — very small changes in $\{\alpha_n\}$ lead to wildly different functions.

In the next section of these notes, we start to get a handle on this process. We start with the concept of an **orthobasis**, where the mapping from vector in a Hilbert space to sequence of coefficients is perfectly stable — changes to the sizes of the coefficients amount to a comparable sized change in reconstructed vector.

Orthogonal bases

A collection of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ in a finite dimensional vector space \mathcal{S} is called an **orthogonal basis** if

1. $\text{span}(\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}) = \mathcal{S}$,
2. $\mathbf{v}_j \perp \mathbf{v}_k$ (i.e. $\langle \mathbf{v}_j, \mathbf{v}_k \rangle = 0$) for all $j \neq k$.

If in addition the vectors are normalized (under the induced norm),

$$\|\mathbf{v}_n\| = 1, \quad \text{for } n = 1, \dots, N,$$

we will call it an **orthonormal basis** or **orthobasis**.

A note on infinite dimensions

In infinite dimensions, we need to be a little more careful with what we mean by “span”. If $\mathcal{B} = \{\mathbf{v}_n\}_{n \in \mathbb{Z}}$ is an infinite sequence of orthogonal vectors in a Hilbert space \mathcal{S} , it is an orthobasis if the *closure* of $\text{span}(\mathcal{B})$ is \mathcal{S} ; this is written

$$\text{cl Span}(\{\mathbf{v}_n\}_n) = \mathcal{S}.$$

We don’t need to get into too much, but basically this means that every vector in \mathcal{S} can be approximated arbitrarily well by a finite linear combination of vectors in \mathcal{B} .

Here is an example which illustrates the point: Let $x(t)$ be any function on $[0, 1]$ which is not a polynomial — say $x(t) = \sin(2\pi t)$. Let $\mathcal{B} = \{1, t, t^2, t^3, \dots\}$; the span (set of a finite linear combinations of elements) of \mathcal{B} is all polynomials on $[0, 1]$. So $\mathbf{x} \notin \text{span}(\mathcal{B})$. But $x(t)$ can be approximated arbitrarily well by elements in \mathcal{B} (using higher and higher order polynomials) so $\mathbf{x} \in \text{cl Span}(\mathcal{B})$.

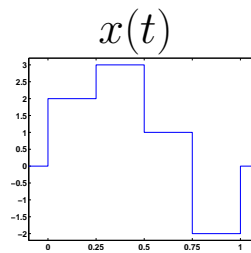
Examples.

1. $\mathcal{S} = \mathbb{R}^2$, equipped with the standard inner product

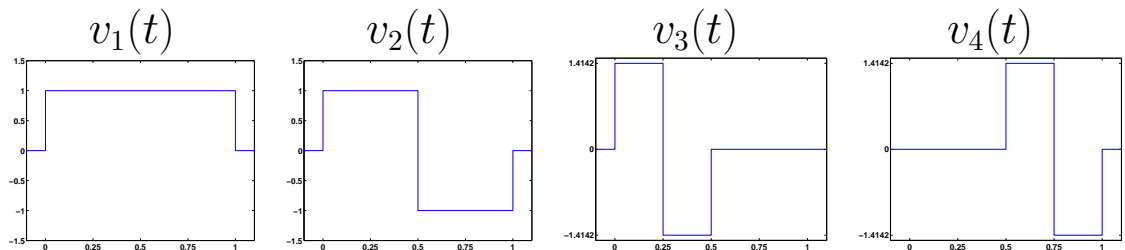
$$\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

2. \mathcal{S} = space of piecewise constant functions on $[0, 1/4)$, $[1/4, 1/2)$, $[1/2, 3/4)$, $[3/4, 1]$

Example signal:



The following four functions form an orthobasis for this space



3. Fourier series

$\{v_k(t) = e^{j2\pi kt}, k \in \mathbb{Z}\}$ is an orthobasis for $L_2([0, 1])$

(with the standard inner product).

Let's quickly check the orthogonality:

$$\begin{aligned}\langle \mathbf{v}_{k_1}, \mathbf{v}_{k_2} \rangle &= \int_0^1 e^{j2\pi(k_1 - k_2)t} dt \\ &= \begin{cases} 1, & k_1 = k_2 \\ 0, & k_1 \neq k_2 \end{cases}.\end{aligned}$$

It is also true that the closure of $\text{span}(\{e^{j2\pi kt}\}_{k=-\infty}^{\infty})$ is $L_2([0, 1])$. The proof of this is a bit involved; if you are interested, see Chapter 5 of Young's *Introduction to Hilbert Space*.

For real-valued functions, we can equivalently use

$$v_0(t) = 1, \quad v_k(t) = \begin{cases} \sqrt{2} \cos(\pi(k+1)t), & k \text{ odd} \\ \sqrt{2} \sin(\pi kt), & k \text{ even} \end{cases}$$

where we have introduced the factors of $\sqrt{2}$ to ensure that $\|\mathbf{v}_k\|_{L_2} = 1$.

4. Legendre Polynomials Define

$$p_0(t) = 1, \quad p_1(t) = t,$$

and then for $n = 1, 2, \dots$

$$p_{n+1}(t) = \frac{2n+1}{n+1} t p_n(t) - \frac{n}{n+1} p_{n-1}(t),$$

and so

$$\begin{aligned} p_2(t) &= \frac{1}{2}(3t^2 - 1) \\ p_3(t) &= \frac{1}{2}(5t^3 - 3t) \\ p_4(t) &= \frac{1}{8}(35t^4 - 30t^2 + 3) \\ &\vdots \quad \text{etc.} \end{aligned}$$

These $p_n(t)$ are called *Legendre polynomials*, and if we renormalize them, taking

$$v_n(t) = \sqrt{\frac{2n+1}{2}} p_n(t),$$

then $v_0(t), \dots, v_N(t)$ are an orthobasis for polynomials of degree N on $[-1, 1]$.

Computing approximations with the Legendre basis is far more stable than computing the approximation in the standard basis.

5. Sampling (Additional reading)

One of the fundamental results in information theory is that a signal that is *bandlimited* can be reconstructed exactly from uniformly spaced samples — this is known as the Shannon-Nyquist sampling theorem, and it has played a major role in the advances in our understanding of data acquisition, imaging, and digital communications.

To state this carefully, we need the notion of a **Fourier transform**. If $f(t)$ is a function of a continuous variable, its Fourier transform is

$$\hat{f}(\omega) = \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt.$$

Roughly speaking, $\hat{f}(\omega)$ is a density that describes how the energy in $f(t)$ is spread over different frequencies. There is a unique mapping between $f(t)$ and $\hat{f}(\omega)$ — given the latter, we can recover $f(t)$ using

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{f}(\omega) e^{j\omega t} d\omega.$$

The mapping also preserves the standard inner product (to within a constant), in that

$$\langle \mathbf{f}, \mathbf{g} \rangle_{L_2} = \int_{-\infty}^{\infty} f(t)g(t) dt = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{f}(\omega)\overline{\hat{g}(\omega)}d\omega = \frac{1}{2\pi} \langle \hat{\mathbf{f}}, \hat{\mathbf{g}} \rangle_{L_2}.$$

(This is known as the classical Parseval identity.)

We say that a function $f(t)$ is **bandlimited** to Ω if

$$\hat{f}(\omega) = 0 \quad \text{for all } |\omega| > \Omega.$$

Roughly speaking, this means that $f(t)$ contains no spectral content at frequencies greater than Ω .

Let $B_\Omega(\mathbb{R}) =$ real-valued functions which are bandlimited to Ω , equipped with the standard inner product. The set of functions

$$\left\{ v_n(t) = \sqrt{T} \frac{\sin(\pi(t - nT)/T)}{\pi(t - nT)}, \quad n \in \mathbb{Z} \right\},$$

with $T = \pi/\Omega$ is an orthobasis for $B_\Omega(\mathbb{R})$. It is a (sort of) easily checked fact that

$$\hat{v}_n(\omega) = \begin{cases} \sqrt{T} e^{-j\omega T n} & |\omega| \leq \Omega \\ 0 & \text{otherwise.} \end{cases}$$

We check the orthogonality of the \mathbf{v}_n :

$$\begin{aligned} & \left\langle \sqrt{T} \frac{\sin(\pi(t - n_1 T)/T)}{\pi(t - n_1 T)}, \sqrt{T} \frac{\sin(\pi(t - n_2 T)/T)}{\pi(t - n_2 T)} \right\rangle \\ &= \frac{1}{2\pi} \int_{-\pi/T}^{\pi/T} T e^{-j\omega T n_1} e^{j\omega T n_2} d\omega \quad (\text{Parseval}) \\ &= \frac{T}{2\pi} \int_{-\pi/T}^{\pi/T} e^{j\omega T (n_1 - n_2)} d\omega \\ &= \begin{cases} 1, & n_1 = n_2 \\ 0, & n_1 \neq n_2 \end{cases}. \end{aligned}$$

That the (closure of the) span of this set is $B_\Omega(\mathbb{R})$ is mathematically equivalent to the Fourier series, where functions on an interval are being composed from harmonic sinusoids. In a few pages, we will show that the expansion coefficients in this basis are samples, thus giving us another way to reinterpret the Shannon-Nyquist sampling theorem.

Linear approximation and orthobases

Let's return to our linear approximation problem:

Given $\mathbf{x} \in \mathcal{S}$, we want to find the closest point in a subspace \mathcal{T} .

Suppose we have an orthobasis $\{\mathbf{v}_1, \dots, \mathbf{v}_N\}$ for \mathcal{T} . Then solving this problem is easy. Here's why: we know the solution is

$$\hat{\mathbf{x}} = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_N \mathbf{v}_N \quad (2)$$

where the α_n are given by

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix} = \mathbf{G}^{-1} \mathbf{b}, \quad \text{with } \mathbf{G} = \begin{bmatrix} \langle \mathbf{v}_1, \mathbf{v}_1 \rangle & \cdots & \langle \mathbf{v}_N, \mathbf{v}_1 \rangle \\ \langle \mathbf{v}_1, \mathbf{v}_2 \rangle & \cdots & \langle \mathbf{v}_N, \mathbf{v}_2 \rangle \\ \vdots & & \vdots \\ \langle \mathbf{v}_1, \mathbf{v}_N \rangle & \cdots & \langle \mathbf{v}_N, \mathbf{v}_N \rangle \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{v}_1 \rangle \\ \langle \mathbf{x}, \mathbf{v}_2 \rangle \\ \vdots \\ \langle \mathbf{x}, \mathbf{v}_N \rangle \end{bmatrix}$$

Now since $\langle \mathbf{v}_n, \mathbf{v}_k \rangle = 1$ if $n = k$ and 0 otherwise, $\mathbf{G} = \mathbf{I}$ (the identity matrix), and so $\mathbf{G}^{-1} = \mathbf{I}$ as well, and

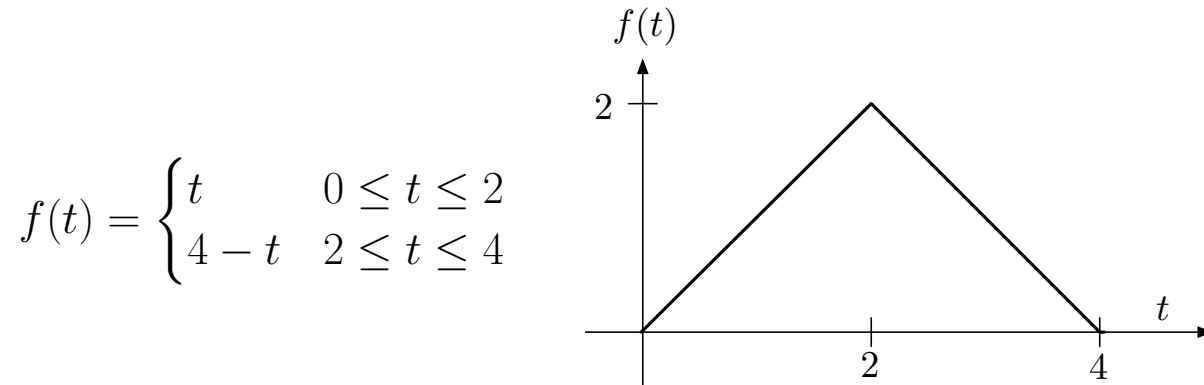
$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{v}_1 \rangle \\ \langle \mathbf{x}, \mathbf{v}_2 \rangle \\ \vdots \\ \langle \mathbf{x}, \mathbf{v}_N \rangle \end{bmatrix}. \quad (3)$$

So calculating the closest point is as easy as computing N inner products — no matrix inversion necessary.

Combining the expressions (2) and (3) gives us the compact expression

$$\hat{\mathbf{x}} = \sum_{n=1}^N \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

Example. Suppose $\mathbf{f} \in L_2([0, 4])$ is



Let \mathcal{T} = piecewise constant functions on $[0, 1)$, $[1, 2)$, $[2, 3)$, $[3, 4]$.

Find the closest point in \mathcal{T} to \mathbf{f} . A good orthobasis to use is

$$v_n(t) = \begin{cases} 1 & (n-1) \leq t \leq n \\ 0 & \text{otherwise} \end{cases}, \quad n = 1, 2, 3, 4.$$

Example Now let

$$f(t) = \begin{cases} t + 1, & -1 \leq t < 0, \\ 1 - t, & 0 \leq t \leq 1. \end{cases}$$

What is the best 4th-order polynomial approximation of \boldsymbol{f} on $[-1, 1]$?

Orthobasis expansions

The orthogonality principle (easily) gives us an expression for the **expansion coefficients** if a vector in an orthobasis.

Suppose a finite dimensional space \mathcal{S} has an orthobasis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$. Given any $\mathbf{x} \in \mathcal{S}$, the closest point in \mathcal{S} to \mathbf{x} is \mathbf{x} itself (of course). This gives us the following **reproducing formula**:

$$\mathbf{x} = \sum_{n=1}^N \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n, \quad \text{for all } \mathbf{x} \in \mathcal{S}.$$

In infinite dimensions, if \mathcal{S} has an orthobasis $\{\mathbf{v}_n\}_{n=-\infty}^{\infty}$ and $\mathbf{x} \in \mathcal{S}$ obeys

$$\sum_{n=-\infty}^{\infty} |\langle \mathbf{x}, \mathbf{v}_n \rangle|^2 < \infty,$$

then we can write

$$\mathbf{x} = \sum_{n=-\infty}^{\infty} \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

(We will write the rigorous version of this shortly.)

In other words, $\mathbf{x} \in \mathcal{S}$ is captured without loss by the discrete list of numbers

$$\dots, \langle \mathbf{x}, \mathbf{v}_{-1} \rangle, \langle \mathbf{x}, \mathbf{v}_0 \rangle, \langle \mathbf{x}, \mathbf{v}_1 \rangle, \dots$$

An orthobasis gives us a natural way to discretize vectors in \mathcal{S} through a set of expansion coefficients. Moreover, there is a straightforward and explicit way to compute these expansion coefficients — you simply take an inner product with the corresponding basis vector.

Example: Change of basis in \mathbb{R}^2 .

Consider \mathbb{R}^2 equipped with the standard Euclidean inner product and the orthonormal basis

$$\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

Let $\mathbf{x} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$. Find α_1, α_2 such that

$$\mathbf{x} = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2.$$

Example: Fourier Series

We have seen that the functions

$$v_k(t) = e^{j2\pi kt},$$

form an orthobasis for $L_2([0, 1])$. Given a function \mathbf{f} on $[0, 1]$, we also know that

$$f(t) = \sum_{k=-\infty}^{\infty} \alpha_k e^{j2\pi kt}, \quad \alpha_k = \int_0^1 f(t) e^{-j2\pi kt} dt.$$

The equation above is equivalent to

$$\mathbf{f} = \sum_{k=-\infty}^{\infty} \langle \mathbf{f}, \mathbf{v}_k \rangle \mathbf{v}_k$$

where $\langle \cdot, \cdot \rangle$ is the L_2 inner product for functions on $[0, 1]$.

For real-valued functions, we can equivalently take

$$f(t) = \sum_{k=0}^{\infty} \alpha_k v_k(t), \quad \alpha_k = \int_0^1 f(t) v_k(t) dt,$$

with

$$v_k(t) = \begin{cases} 1, & k = 0, \\ \sqrt{2} \cos(\pi(k+1)t), & k \text{ odd}, \\ \sqrt{2} \sin(\pi kt), & k \text{ even}. \end{cases}$$

Example: Change of basis for a polynomial

Let $f(t) = 5t^2 - 3t + 2$ for $t \in [0, 1]$. Find $\alpha_0, \alpha_1, \alpha_2$ such that

$$f(t) = \alpha_0 v_0(t) + \alpha_1 v_1(t) + \alpha_2 v_2(t),$$

where the $v_n(t)$ are the Legendre polynomials given earlier in these notes.

Example: Sampling a bandlimited function (additional reading).

$B_{\pi/T}$ = space of bandlimited signals equipped with the standard inner product. We have seen already that

$$v_n(t) = \sqrt{T} \frac{\sin(\pi(t - nT)/T)}{\pi(t - nT)}, \quad n \in \mathbb{Z}$$

is an orthobasis for $B_{\pi/T}$. This means that any $\mathbf{f} \in B_{\pi/T}$ can be written

$$\mathbf{f} = \sum_{n=-\infty}^{\infty} \langle \mathbf{f}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

What are the $\langle \mathbf{f}, \mathbf{v}_n \rangle$? As we saw above

$$\begin{aligned} \langle \mathbf{f}, \mathbf{v}_n \rangle &= \left\langle f(t), \sqrt{T} \frac{\sin(\pi(t - nT)/T)}{\pi(t - nT)} \right\rangle \\ &= \frac{1}{2\pi} \int_{-\pi/T}^{\pi/T} \hat{f}(\omega) \sqrt{T} e^{jn\omega T} d\omega \\ &= \sqrt{T} f(nT), \end{aligned}$$

which is simply a sample scaled by \sqrt{T} . The reproducing formula in this case is a restatement of the Shannon-Nyquist sampling theorem:

$$\begin{aligned} f(t) &= \sum_{n=-\infty}^{\infty} \langle \mathbf{f}, \mathbf{v}_n \rangle \mathbf{v}_n \\ &= \sum_{n=-\infty}^{\infty} f(nT) \frac{T \sin(\pi(t - nT)/T)}{\pi(t - nT)}, \end{aligned}$$

meaning that knowing $f(t)$ only at the discrete set of points $\{nT, n \in \mathbb{Z}\}$ is enough to recover $f(t)$ for all $t \in \mathbb{R}$ using “sinc interpolation”.

The moral of the story is that we can recreate a vector \mathbf{x} in a Hilbert space from the sequence of numbers $\{\langle \mathbf{x}, \mathbf{v}_n \rangle\}$. We can think of every different orthobasis for \mathcal{S} as a different **transform**, and the $\{\langle \mathbf{x}, \mathbf{v}_n \rangle\}$ as **transform coefficients**. Let us put down the technical details under which this is true.

Technical details: Convergence of orthonormal sequences in Hilbert space

We now solidify the mathematics to make sense of what we have called the reproducing formula

$$\mathbf{x} = \sum_{n=1}^{\infty} \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

There are an infinite number of terms in this sum, so we will have to take care in showing that the sum converges to something meaningful. Let \mathbf{x}_N be the partial sum

$$\mathbf{x}_N = \sum_{n=1}^N \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

We want to show that the sequence

$$\{\mathbf{x}_N\}_{N=1}^{\infty} = \{\mathbf{x}_1, \mathbf{x}_2, \dots\}$$

converges to \mathbf{x} under the norm induced by the inner product. Technically, this means that for every $\epsilon > 0$, there exists an m such that

$$\|\mathbf{x} - \mathbf{x}_N\| \leq \epsilon \quad \text{for all } N \geq m.$$

We usually write this as

$$\lim_{N \rightarrow \infty} \|\mathbf{x} - \mathbf{x}_N\| = 0$$

or $\mathbf{x}_N \rightarrow \mathbf{x}$ as $N \rightarrow \infty$.

We establish this with a series of three propositions. The first two are technical, and are used as an end to prove the third.

Proposition 1: Inner product with a fixed vector is a continuous function. Let \mathbf{c} be a vector in a Hilbert space \mathcal{S} , and let $\{\mathbf{x}_N\}$ be a sequence of vectors in \mathcal{S} that converge to a point \mathbf{x} . Then

$$\lim_{N \rightarrow \infty} \langle \mathbf{x}_N, \mathbf{c} \rangle = \langle \mathbf{x}, \mathbf{c} \rangle.$$

Proof. This follows directly from Cauchy-Schwarz:

$$\begin{aligned} |\langle \mathbf{x}_N, \mathbf{c} \rangle - \langle \mathbf{x}, \mathbf{c} \rangle| &= |\langle \mathbf{x}_N - \mathbf{x}, \mathbf{c} \rangle| \\ &\leq \|\mathbf{x}_N - \mathbf{x}\| \|\mathbf{c}\|. \end{aligned}$$

Thus

$$\|\mathbf{x}_N - \mathbf{x}\| \leq \frac{\epsilon}{\|\mathbf{c}\|} \quad \Rightarrow \quad |\langle \mathbf{x}_N, \mathbf{c} \rangle - \langle \mathbf{x}, \mathbf{c} \rangle| \leq \epsilon,$$

and so $\langle \mathbf{x}_N, \mathbf{c} \rangle \rightarrow \langle \mathbf{x}, \mathbf{c} \rangle$ as $\mathbf{x}_N \rightarrow \mathbf{x}$. ■

Proposition 2: Bessel's inequality. Let $\{\mathbf{v}_n\}_{n=1}^{\infty}$ be a sequence of orthonormal vectors in a Hilbert space \mathcal{S} . Then for any $\mathbf{x} \in \mathcal{S}$,

$$\sum_{n=1}^{\infty} |\langle \mathbf{x}, \mathbf{v}_n \rangle|^2 \leq \|\mathbf{x}\|^2. \quad (4)$$

Proof. Fix N and let

$$\mathbf{y}_N = \sum_{n=1}^N \langle \mathbf{x}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

We know that \mathbf{y}_N is the closest point in $\text{Span}(\{\mathbf{v}_1, \dots, \mathbf{v}_N\})$ to \mathbf{x} , and that

$$\langle \mathbf{x} - \mathbf{y}_N, \mathbf{y}_N \rangle = 0.$$

Thus by the Pythagorean rule,

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

Since $\|\mathbf{y}_N\|^2 = \sum_{n=1}^N |\langle \mathbf{x}, \mathbf{v}_n \rangle|^2$, we have

$$\begin{aligned} \sum_{n=1}^N |\langle \mathbf{x}, \mathbf{v}_n \rangle|^2 &= \|\mathbf{x}\|^2 - \|\mathbf{x} - \mathbf{y}_N\|^2 \\ &\leq \|\mathbf{x}\|^2. \end{aligned}$$

Since the bound holds uniformly for all N , (4) follows. ■

Note that we are not requiring $\{\mathbf{v}_n\}_{n=1}^\infty$ to span the entire space — if it does span the space, it is an orthonormal basis, and the inequality in (4) becomes equality.

Proposition 3: The reproducing formula. Let $\{\mathbf{v}_n\}_{n=1}^\infty$ be a sequence of orthonormal vectors in a Hilbert space \mathcal{S} , and let $\{\alpha_n\}$ be a sequence of scalars. Set

$$\mathbf{x}_N = \sum_{n=1}^N \alpha_n \mathbf{v}_n. \tag{5}$$

Then the sequence of vectors $\{\mathbf{x}_N\}$ converges to a vector $\mathbf{x} \in \mathcal{S}$ as $N \rightarrow \infty$ if and only if

$$\sum_{n=1}^{\infty} |\alpha_n|^2 < \infty. \tag{6}$$

Moreover, the point \mathbf{x} obeys

$$\langle \mathbf{x}, \mathbf{v}_n \rangle = \alpha_n \quad \text{for all } n = 1, 2, \dots$$

Proof. By “converge” we mean that there is a vector $\mathbf{x} \in \mathcal{S}$ such that

$$\lim_{N \rightarrow \infty} \|\mathbf{x} - \mathbf{x}_N\|_{\mathcal{S}} = 0.$$

First, suppose that (6) holds, and take \mathbf{x}_N as in (5). Then for any $m \geq 1$,

$$\begin{aligned} \|\mathbf{x}_{N+m} - \mathbf{x}_N\|^2 &= \left\| \sum_{n=N+1}^{N+m} \alpha_n \mathbf{v}_n \right\|^2 \\ &= \sum_{n=N+1}^{N+m} |\alpha_n|^2, \end{aligned}$$

where the second step follows from the Pythagorean theorem and the fact that $\|\mathbf{v}_n\| = 1$. Since the $\{\alpha_n\}$ are square-summable, we have

$$\|\mathbf{x}_{N+m} - \mathbf{x}_N\|^2 \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

Thus $\{\mathbf{x}_N\}_{N=1}^{\infty}$ is a Cauchy sequence, and since \mathcal{S} is a Hilbert space, $\{\mathbf{x}_N\}$ converges.

Now suppose that the sequence $\{\mathbf{x}_N\}$ in (5) converges; that is, there exists a \mathbf{x} such that

$$\lim_{N \rightarrow \infty} \|\mathbf{x} - \mathbf{x}_N\| = 0.$$

For a fixed $k \leq N$, we have

$$\begin{aligned} \langle \mathbf{x}_N, \mathbf{v}_k \rangle &= \sum_{n=1}^N \alpha_n \langle \mathbf{v}_n, \mathbf{v}_k \rangle \\ &= \alpha_k, \end{aligned}$$

since $\langle \mathbf{v}_n, \mathbf{v}_k \rangle = 1$ when $n = k$ and is zero otherwise. Since the above holds for all $N \geq k$, we have

$$\begin{aligned}\alpha_k &= \lim_{N \rightarrow \infty} \langle \mathbf{x}_N, \mathbf{v}_k \rangle \\ &= \left\langle \lim_{N \rightarrow \infty} \mathbf{x}_N, \mathbf{v}_k \right\rangle \\ &= \langle \mathbf{x}, \mathbf{v}_k \rangle.\end{aligned}$$

(The first step above follows from the continuity on the inner product; see the lemma above.) Then applying the Bessel inequality,

$$\begin{aligned}\sum_{n=1}^{\infty} |\alpha_n|^2 &= \sum_{n=1}^{\infty} |\langle \mathbf{x}, \mathbf{v}_n \rangle|^2 \\ &\leq \|\mathbf{x}\|^2.\end{aligned}$$

Since we started with the supposition that $\mathbf{x} \in \mathcal{S}$, we know $\|\mathbf{x}\| \leq \infty$.

A useful property: Parseval's Theorem

Once an orthobasis is introduced into any Hilbert space, it can be mapped to standard Euclidean space. As we will see, inner products between vectors become standard dot products between their expansion coefficients, and (induced) norms become standard sum-of-squares norms. The key identity is called Parseval's theorem¹, and it is relatively easy to establish.

Let \mathcal{S} be a Hilbert space with inner product $\langle \cdot, \cdot \rangle_{\mathcal{S}}$ which induces the norm $\| \cdot \|_{\mathcal{S}}$. Let $\{v_k\}_{k \in \Gamma}$ be an orthobasis² for \mathcal{S} . Then for every $\mathbf{x}, \mathbf{y} \in \mathcal{S}$,

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{S}} = \sum_{k \in \Gamma} \alpha_k \overline{\beta_k},$$

where

$$\alpha_k = \langle \mathbf{x}, \mathbf{v}_k \rangle_{\mathcal{S}}, \quad \beta_k = \langle \mathbf{y}, \mathbf{v}_k \rangle_{\mathcal{S}}.$$

You can think of the $\{\alpha_k\}$ as the transform coefficients of \mathbf{x} and the $\{\beta_k\}$ as the transform coefficients of \mathbf{y} . So we have

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{S}} &= \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle_{\ell_2}, \\ \|\mathbf{x}\|_{\mathcal{S}}^2 &= \|\boldsymbol{\alpha}\|_2^2. \end{aligned}$$

\Rightarrow An orthobasis makes every Hilbert space **equivalent** to ℓ_2 .

¹In its most common usage, the name Parseval is usually associated with the preservation of energy (to within a constant) between a function $f(t)$ and its Fourier transform, as we discussed in the last set of notes. Here, we will use this terminology more broadly.

²We are using Γ to be an arbitrary index set here; it can be either finite, e.g. $\Gamma = 1, 2, \dots, N$, or infinite, e.g. $\Gamma = \mathbb{Z}$.

All of the geometry (lengths, angles) maps into standard Euclidean geometry in coefficient space. As you can imagine, this is a pretty useful fact.

Proof of Parseval. With $\alpha_k = \langle \mathbf{x}, \mathbf{v}_k \rangle$ and $\beta_k = \langle \mathbf{y}, \mathbf{v}_k \rangle$, we can write

$$\mathbf{x} = \sum_{k \in \Gamma} \alpha_k \mathbf{v}_k, \quad \text{and} \quad \mathbf{y} = \sum_{k \in \Gamma} \beta_k \mathbf{v}_k,$$

and so

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} \rangle_S &= \left\langle \sum_{k \in \Gamma} \alpha_k \mathbf{v}_k, \sum_{\ell \in \Gamma} \beta_\ell \mathbf{v}_\ell \right\rangle_S \\ &= \sum_{k \in \Gamma} \alpha_k \left\langle \mathbf{v}_k, \sum_{\ell \in \Gamma} \beta_\ell \mathbf{v}_\ell \right\rangle_S \\ &= \sum_{k \in \Gamma} \sum_{\ell \in \Gamma} \alpha_k \bar{\beta}_\ell \langle \mathbf{v}_k, \mathbf{v}_\ell \rangle_S. \end{aligned}$$

For a fixed value of k , only one term in the inner sum above will be nonzero, as $\langle \mathbf{v}_k, \mathbf{v}_\ell \rangle = 0$ unless $\ell = k$. Thus

$$\langle \mathbf{x}, \mathbf{y} \rangle_S = \sum_{k \in \Gamma} \alpha_k \bar{\beta}_k.$$

A straightforward consequence of the result above is that distances in \mathcal{S} under the induced norm are equivalent to Euclidean (ℓ_2) distances in coefficient space:

$$\|\mathbf{x} - \mathbf{y}\|_S = \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|_2 = \left(\sum_{k \in \Gamma} (\alpha_k - \beta_k)^2 \right)^{1/2}.$$

Thus changing the value of an orthobasis expansion coefficient by an amount ϵ will change the signal by an amount (as measured in $\|\cdot\|_S$) ϵ .

To be more precise about this, suppose \mathbf{x} has transform coefficients $\{\alpha_k = \langle \mathbf{x}, \mathbf{v}_k \rangle_S\}$. If I perturb one of them, say at location k_0 , by setting

$$\tilde{\alpha}_k = \begin{cases} \alpha_{k_0} + \epsilon & k = k_0 \\ \alpha_k & k \neq k_0 \end{cases},$$

and then synthesizing

$$\tilde{\mathbf{x}} = \sum_{k \in \Gamma} \tilde{\alpha}_k \mathbf{v}_k,$$

we will have

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_S = \epsilon.$$

Notice that while the error is localized to one expansion coefficient, it could effect the entire reconstruction, but its net effect will still be ϵ .

More generally, suppose we add an error to every coefficient:

$$\tilde{\alpha}_k = \alpha_k + \epsilon_k,$$

then re-synthesize the vector

$$\tilde{\mathbf{x}} = \sum_{k \in \Gamma} \tilde{\alpha}_k \mathbf{v}_k.$$

The total error (as measured by the norm $\|\cdot\|_{\mathcal{S}}$) in the reconstructed vector will be the same as the total error (as measured by the Euclidean norm $\|\cdot\|_2$):

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_{\mathcal{S}} = \|\boldsymbol{\alpha} - \tilde{\boldsymbol{\alpha}}\|_2 = \left(\sum_{k \in \Gamma} |\epsilon_k|^2 \right)^{1/2}.$$

The upshot of this is that as we manipulate the expansion coefficients $\{\alpha_k\}$, we know what the net effect will be in the reconstructed function.

Truncating ortho expansions and linear approximation

Say $\{\mathbf{v}_k\}_{k=0}^{\infty}$ is an orthobasis for a Hilbert space \mathcal{S} . Let \mathcal{T} be the subspace spanned by the first 10 elements of $\{\mathbf{v}_k\}$:

$$\mathcal{T} = \text{span}(\{\mathbf{v}_0, \dots, \mathbf{v}_9\}).$$

1. Given $\mathbf{x} \in \mathcal{S}$, what is the closest point in \mathcal{T} (call it $\hat{\mathbf{x}}$) to \mathbf{x} ?
We have seen that it is

$$\hat{\mathbf{x}} = \sum_{k=0}^9 \langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k,$$

where $\langle \cdot, \cdot \rangle$ is the inner product for \mathcal{S} .

2. How good an approximation is $\hat{\mathbf{x}}$ to \mathbf{x} ? If we measure this in the induced norm $\|\cdot\|$, then

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|^2 &= \left\| \sum_{k=0}^{\infty} \langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k - \sum_{k=0}^9 \langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k \right\|^2 \\ &= \left\| \sum_{k=10}^{\infty} \langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k \right\|^2 \\ &= \sum_{k=10}^{\infty} |\langle \mathbf{x}, \mathbf{v}_k \rangle|^2.\end{aligned}$$

Since also

$$\|\mathbf{x}\|^2 = \sum_{k=0}^{\infty} |\langle \mathbf{x}, \mathbf{v}_k \rangle|^2$$

the approximation error for $\hat{\mathbf{x}}$ will be small if the first 10 transform coefficients

$$\langle \mathbf{x}, \mathbf{v}_0 \rangle, \langle \mathbf{x}, \mathbf{v}_1 \rangle, \dots, \langle \mathbf{x}, \mathbf{v}_9 \rangle,$$

contain “most” of the total energy.

Of course, there is nothing special about taking the first 10 coefficients. We can just as easily form a K term approximation using

$$\hat{\mathbf{x}}_K = \sum_{k=0}^{K-1} \langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k$$

which has error

$$\|\mathbf{x} - \hat{\mathbf{x}}_K\|^2 = \sum_{k=K}^{\infty} |\langle \mathbf{x}, \mathbf{v}_k \rangle|^2.$$

If the sum above is small for moderately large K , we can “compress” \mathbf{x} by using just the first K terms in the expansion.

This is precisely what is done in image and video compression — more details below.

Example:

Any real-valued function on $[-1/2, 1/2]$ with even symmetry can be built up out of harmonic cosines:

$$f(t) = \alpha_0 + \sum_{k=1}^{\infty} \alpha_k \sqrt{2} \cos(2\pi kt).$$

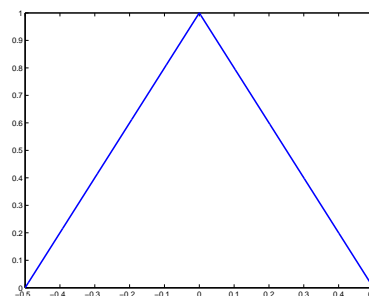
(That this is true follows directly from the observation that every signal on $[-1/2, 1/2]$ that is real-valued and even has a Fourier series which is real-valued and even.) This is an orthobasis expansion in the standard inner product with

$$v_0(t) = 1, \quad v_1(t) = \sqrt{2} \cos(2\pi t), \quad \dots, \quad v_k(t) = \sqrt{2} \cos(2\pi kt), \quad \dots$$

It is easy to check that $\langle \mathbf{v}_k, \mathbf{v}_\ell \rangle = 0$, $k \neq \ell$ and $\langle \mathbf{v}_k, \mathbf{v}_k \rangle = 1$.

For the triangle function below

$$f(t) = \begin{cases} 1 + 2t, & -1/2 \leq t \leq 0 \\ 1 - 2t, & 0 \leq t \leq 1/2 \end{cases}$$



the expansion coefficients are

$$\begin{aligned}
 \alpha_0 &= 1/2, \\
 \alpha_k &= \int_{-1/2}^{1/2} f(t) \sqrt{2} \cos(2\pi kt) dt \\
 &= 2\sqrt{2} \int_0^{1/2} (1-2t) \cos(2\pi kt) dt \\
 &= \begin{cases} 0 & k \text{ even, } k \neq 0 \\ \frac{2\sqrt{2}}{\pi^2 k^2} & k \text{ odd} \end{cases}.
 \end{aligned}$$

First, let's compute the norm in time and coefficient space just to make sure they agree:

$$\|\mathbf{f}\|_2^2 = \int_{-1/2}^{1/2} |f(t)|^2 dt = 2 \int_0^{1/2} (1-2t)^2 dt = 1/3,$$

and

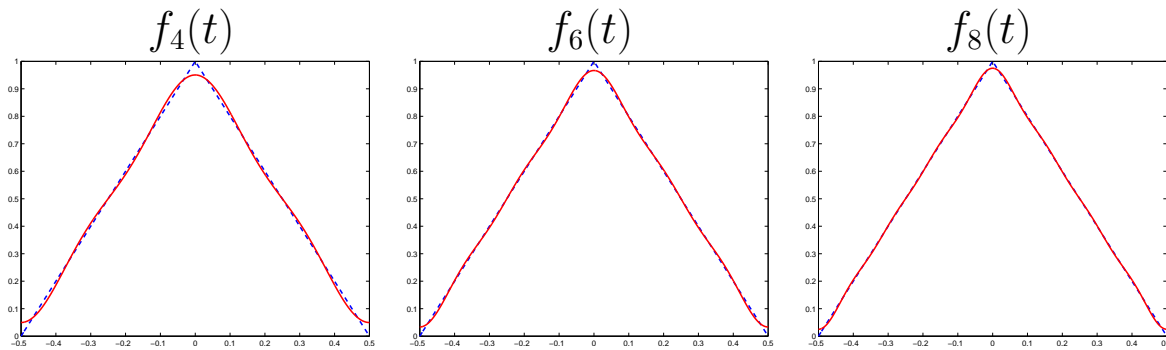
$$\begin{aligned}
 \sum_{k=0}^{\infty} |\alpha_k|^2 &= \frac{1}{4} + \frac{8}{\pi^4} \sum_{k'=0}^{\infty} \frac{1}{(1+2k')^4} \\
 &= \frac{1}{4} + \frac{8}{\pi^4} \left(\frac{\pi^4}{96} \right) \\
 &= \frac{1}{3}.
 \end{aligned}$$

When we truncate the expansion at K terms,

$$f_K(t) = \frac{1}{2} + \sum_{k=1}^{K-1} \alpha_k \sqrt{2} \cos(2\pi kt),$$

we can interpret the result as an **approximation** of $f(t)$ that is a member of the K -dimensional subspace $\text{span}(\{\sqrt{2}\cos(2\pi kt)\}_{k=0}^{K-1})$, and we know that it is the best approximation in that subspace.

Here are the approximations for $K = 4, 6, 8$:



We can compute the error in each of these approximations explicitly, as

$$\begin{aligned} f(t) - f_K(t) &= \sum_{k=0}^{\infty} \alpha_k \sqrt{2} \cos(2\pi kt) - \sum_{k=0}^{K-1} \alpha_k \sqrt{2} \cos(2\pi kt) \\ &= \sum_{k=K}^{\infty} \alpha_k \sqrt{2} \cos(2\pi kt), \end{aligned}$$

and so

$$\|\mathbf{f} - \mathbf{f}_K\|_2^2 = \sum_{k=K}^{\infty} |\alpha_k|^2,$$

or, since $\mathbf{f}_K \perp \mathbf{f} - \mathbf{f}_K$,

$$\|\mathbf{f} - \mathbf{f}_K\|_2^2 = \|\mathbf{f}\|_2^2 - \|\mathbf{f}_K\|_2^2.$$

In the three examples above, we have

$$\begin{aligned} \|\mathbf{f} - \mathbf{f}_4\|_2^2 &\approx 1.92 \cdot 10^{-4}, \quad \|\mathbf{f} - \mathbf{f}_6\|_2^2 \approx 6.01 \cdot 10^{-5}, \\ \|\mathbf{f} - \mathbf{f}_8\|_2^2 &\approx 2.59 \cdot 10^{-5}. \end{aligned}$$

Application: The DCT and JPEG

A great example of how manipulating orthobasis expansion coefficients can lead to something useful is given by the JPEG image compression standard. This is a complicated compression standard, but it is based on a simple idea: break the image into pieces, represent each piece using a cosine basis, then achieve compression by truncating/quantizing the expansion.

Lets start by recalling the discrete Fourier transform for vectors in \mathbb{C}^N . This is basically the same as Fourier series, except that the sinusoidal basis vectors are discrete (instead of functions of a continuous variable). Any $\mathbf{x} \in \mathbb{C}^N$ can be written as

$$\mathbf{x} = \sum_{k=0}^{N-1} \alpha_k \boldsymbol{\psi}_k, \quad \psi_k[n] = \frac{1}{\sqrt{N}} e^{j2\pi nk/N}, \quad n = 0, \dots, N-1.$$

Computing $\boldsymbol{\alpha}$ for a given \mathbf{x} can be done in $O(N \log_2 N)$ time using the *fast Fourier transform* (FFT), one of the most important algorithms in engineering and applied mathematics.

Closely related is the *discrete cosine transform* (DCT). This is a basis for \mathbb{R}^N , and has slightly more favorable symmetry properties than the standard DFT. The DCT basis functions for \mathbb{R}^N are

$$\psi_k[n] = \begin{cases} \sqrt{\frac{1}{N}} & k = 0 \\ \sqrt{\frac{2}{N}} \cos\left(\frac{\pi k}{N}\left(n + \frac{1}{2}\right)\right) & k = 1, \dots, N-1 \end{cases} \quad (7)$$

for sample indices $n = 0, 1, \dots, N-1$. Showing that

$$\sum_{n=0}^{N-1} \psi_k[n] \psi_\ell[n] = \begin{cases} 1 & k = \ell \\ 0 & k \neq \ell \end{cases}$$

is an exercise you can do at home. Notice that the samples of the cosines are on the half-sample points (we see $(n + 1/2)$ in the expression above instead of n).

The DCT can be quickly computed from the DFT of a symmetric extension of \mathbf{x} . That means we have a **fast algorithm** for computing the DCT — the cost is essentially the same as for an FFT, $O(N \log N)$.

An *image* is an array of pixels arranged on a grid indexed by n_1, n_2 . We can think of $x[n_1, n_2]$ as signifying the intensity of an image at location (n_1, n_2) . If we want more than black and white images, we can use three intensity arrays (three “channels”) to represent the color at every point.

Here, we will assume that these arrays are square, $N \times N$, although the discussion below is easily adapted for rectangular arrays. We will still call $\{x[n_1, n_2], 0 \leq n_1, n_2 \leq N - 1\}$ a vector, even though it is naturally indexed by two variables. (It should be clear anyway that arrays of numbers like this obey our requirements for an abstract vector space.) We will call this space of 2D arrays $\mathbb{R}^N \times \mathbb{R}^N$.

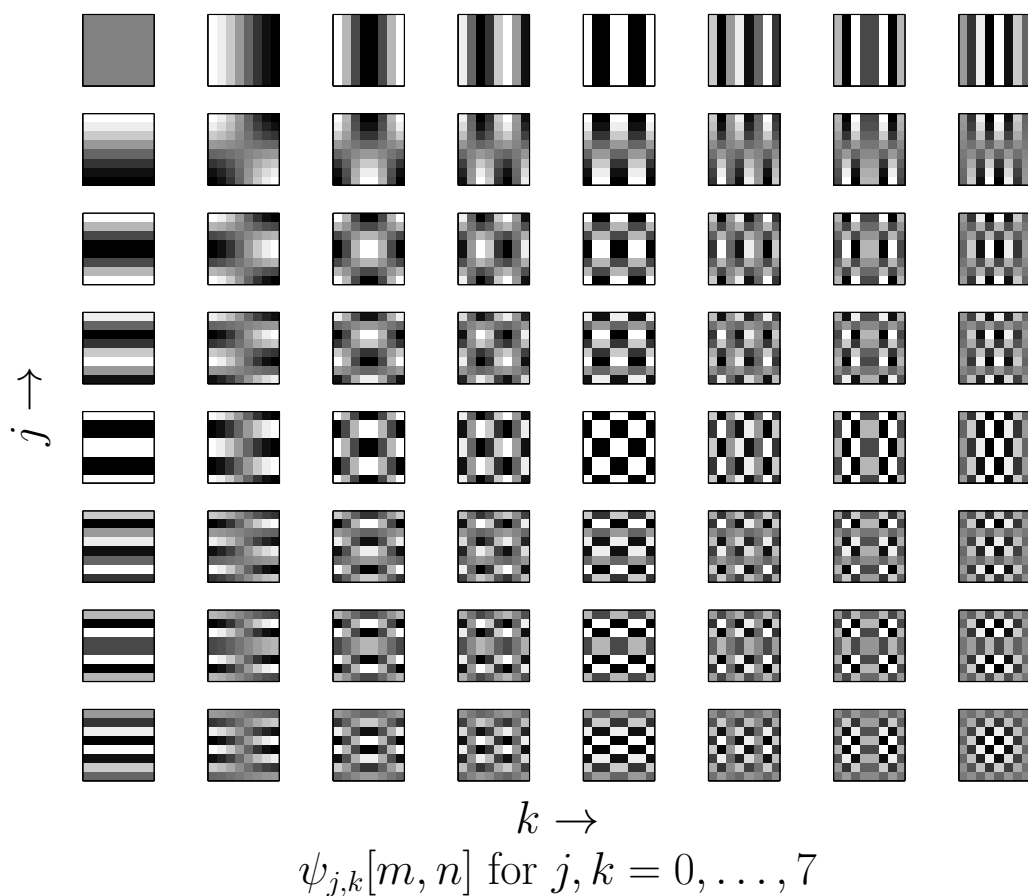
There is a straightforward way to create an orthobasis for $\mathbb{R}^N \times \mathbb{R}^N$ from an orthobasis for \mathbb{R}^N . We simply create *separable 2D arrays* from the 1D basis vectors. More precisely, if $\{\boldsymbol{\psi}_k, k = 0, \dots, N - 1\}$ is an orthobasis for \mathbb{R}^N , then $\{\boldsymbol{\psi}_{j,k}^{2D}, 0 \leq j, k \leq N\}$ is an orthobasis for $\mathbb{R}^N \times \mathbb{R}^N$, where

$$\psi_{j,k}^{2D}[n_1, n_2] = \psi_j[n_1]\psi_k[n_2].$$

You will prove this (and actually something much more general than

this) on next week's homework.

Let's take a look at the 2D DCT orthobasis for 8×8 image patches. The 64 2D DCT basis functions for $N = 8$ are shown below:



2D DCT coefficients are indexed by two integers, and so are naturally arranged on a grid as well:

$$\begin{array}{cccc}
 \alpha_{0,0} & \alpha_{0,1} & \cdots & \alpha_{0,N-1} \\
 \alpha_{1,0} & \alpha_{1,1} & \cdots & \alpha_{1,N-1} \\
 \vdots & \vdots & \vdots & \vdots \\
 \alpha_{N-1,0} & \alpha_{N-1,1} & \cdots & \alpha_{N-1,N-1}
 \end{array}$$

The DCT in image and video compression

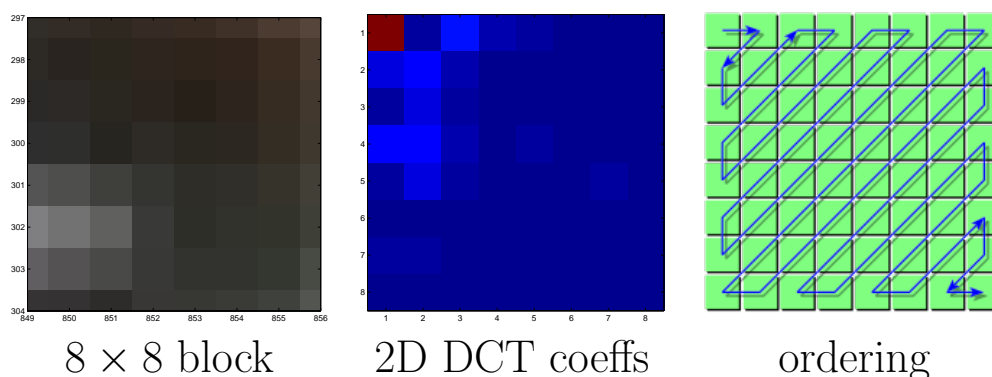
The DCT is basis of the popular JPEG image compression standard. The central idea is that while energy in a picture is distributed more or less evenly throughout, in the DCT transform domain it tends to be *concentrated* at low frequencies.

JPEG compression work roughly as follows:

1. Divide the image into 8×8 blocks of pixels
2. Take a DCT within each block
3. Quantize the coefficients — the rough effect of this is to keep the larger coefficients and remove the smaller ones
4. Bitstream (losslessly) encode the result.

There are some details we are leaving out here, probably the most important of which is how the three different color bands are dealt with, but the above outlines the essential ideas.

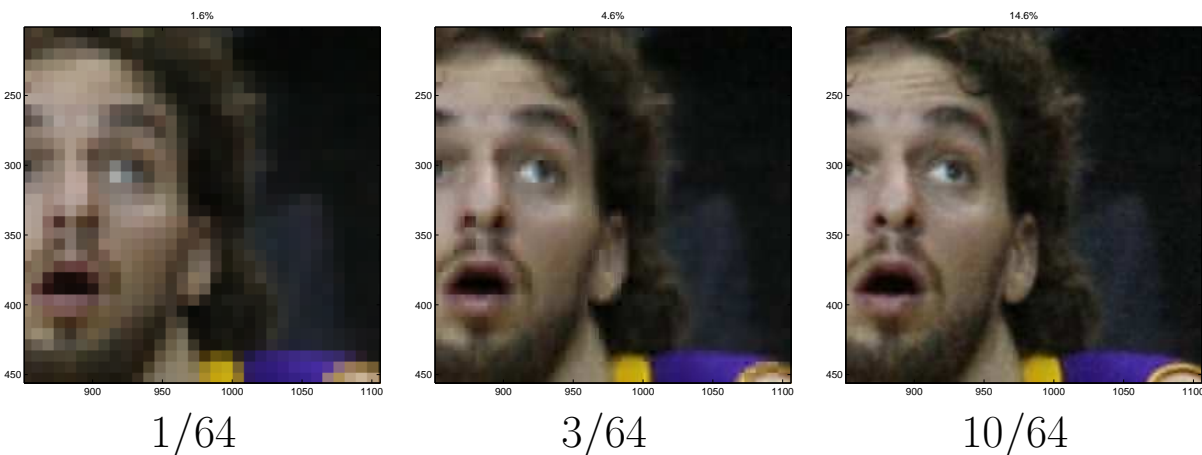
The basic idea is that while the energy within an 8×8 block of pixels tends to be more or less evenly distributed, the DCT concentrates this energy onto a relatively small number of transform coefficients. Moreover, the significant coefficients tend to be at the same place in the transform domain (low spatial frequencies).



To get a rough feel for how closely this model matches reality, let's look at a simple example. Here we have an original image 2048×2048 , and a zoom into a 256×256 piece of the image:



Here is the same piece after using 1 of the 64 coefficients per block ($1/64 \approx 1.6\%$), $3/64 \approx 4.6\%$ of the coefficients, and $10/64 \approx 15.62\%$:



So the “low frequency” heuristic appears to be a good one.

JPEG does not just “keep or kill” coefficients in this manner, it quantizes them using a fixed quantization mask. Here is a common

example:

$$Q = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}.$$

The quantization simply maps $\alpha_{j,k} \rightarrow \tilde{\alpha}_{j,k}$ using

$$\tilde{\alpha}_{j,k} = Q_{j,k} \cdot \text{round} \left(\frac{\alpha_{j,k}}{Q_{j,k}} \right)$$

You can see that the coefficients at low frequencies (upper left) are being treated much more gently than those at higher frequencies (lower right).

The **decoder** simply reconstructs each 8×8 block \mathbf{x}_b using the synthesis formula

$$\tilde{\mathbf{x}}_b[m, n] = \sum_{k=0}^7 \sum_{\ell=0}^7 \tilde{\alpha}_{k,\ell} \phi_{k,\ell}[m, n]$$

By the Parseval theorem, we know exactly what the effect of quantizing each coefficient is going to be on the total error, as

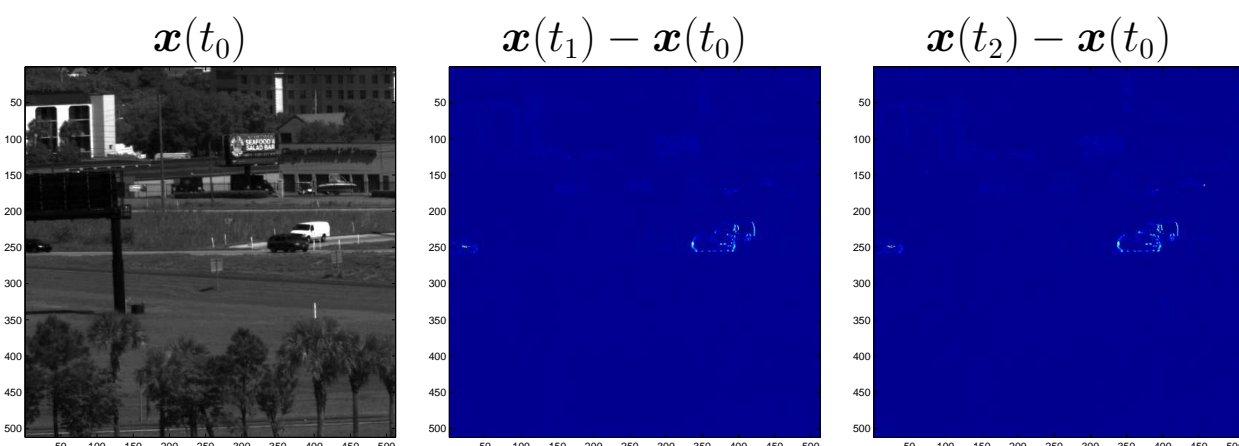
$$\|\mathbf{x}_b - \tilde{\mathbf{x}}_b\|_2^2 = \|\boldsymbol{\alpha} - \tilde{\boldsymbol{\alpha}}\|_2^2 = \sum_{k=0}^7 \sum_{\ell=0}^7 |\alpha_{k,\ell} - \tilde{\alpha}_{k,\ell}|^2.$$

Video compression

The DCT also plays a fundamental role in video compression (e.g. MPEG, H.264, etc.), but in a slightly different way. Video codecs are complicated, but here is essentially what they do:

1. Estimate, describe, and quantize the motion in between frames.
2. Use the motion estimate to “predict” the next frame.
3. Use the (block-based) DCT to code the residual.

Here is an example video frame, along with the differences between this frame and the next two frames (in false color):



The only activity is where the car is moving from left to right.