# Contents

# 3   Inner Product Spaces

In this chapter we will study vector spaces having an additional kind of structure called an inner product, which generalizes the idea of the dot product of vectors in $\mathbb{R}^n$, and which will allow us to formulate notions of "length" and "angle" in more general vector spaces.

We define inner products in real and complex vector spaces and establish some of their properties, including the celebrated Cauchy-Schwarz inequality, and survey some applications.

We then discuss orthogonality of vectors and subspaces and in particular describe a method for constructing an orthonormal basis for any finite-dimensional inner product space, which provides an analogue of giving "standard unit coordinate axes" in $\mathbb{R}^n$.

Next, we discuss some interactions between inner products and linear transformations, and introduce the fundamental notion of the adjoint of a linear transformation.

Finally, we close with a discussion of two very important practical applications of inner products and orthogonality: computing least-squares approximations and approximating periodic functions with Fourier series.

## 3.1   Inner Product Spaces

- <u>Notation</u>: In this chapter, we will use parentheses around vectors rather than angle brackets, since we will shortly be using angle brackets to denote an inner product. We will also avoid using dots when discussing scalar multiplication, and reserve the dot notation for the dot product of two vectors.

### 3.1.1 Inner Products on Real Vector Spaces

- Our first goal is to describe the natural analogue in an arbitrary real vector space of the dot product in $\mathbb{R}^n$. We recall some basic properties of the dot product:

    ○ The dot product distributes over addition and scaling: $(\mathbf{v}_1 + c\mathbf{v}_2) \cdot \mathbf{w} = (\mathbf{v}_1 \cdot \mathbf{w}) + c(\mathbf{v}_2 \cdot \mathbf{w})$ for any vectors $\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}$ and scalar $c$.

    ○ The dot product is commutative: $\mathbf{v} \cdot \mathbf{w} = \mathbf{w} \cdot \mathbf{v}$ for any vectors $\mathbf{v}, \mathbf{w}$.

    ○ The dot product is nonnegative: $\mathbf{v} \cdot \mathbf{v} \geq 0$ for any vector $\mathbf{v}$.

- We will use these three properties to define an arbitrary inner product on a real vector space:

- Definition: If $V$ is a real vector space, an inner product on $V$ is a pairing that assigns a scalar in $F$ to each ordered pair $(\mathbf{v}, \mathbf{w})$ of vectors in $V$. This pairing is denoted $\langle \mathbf{v}, \mathbf{w} \rangle$ and must satisfy the following properties:

    **[I1]** Linearity in the first argument: $\langle \mathbf{v}_1 + c\mathbf{v}_2, \mathbf{w} \rangle = \langle \mathbf{v}_1, \mathbf{w} \rangle + c\langle \mathbf{v}_2, \mathbf{w} \rangle$.

    **[I2]** Symmetry: $\langle \mathbf{w}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$.

    **[I3]** Positive-definiteness: $\langle \mathbf{v}, \mathbf{v} \rangle \geq 0$ for all $\mathbf{v}$, and $\langle \mathbf{v}, \mathbf{v} \rangle = 0$ only when $\mathbf{v} = \mathbf{0}$.

    ○ The linearity and symmetry properties are fairly clear: if we fix the second component, the inner product behaves like a linear function in the first component, and we want both components to behave in the same way.

    ○ The positive-definiteness property is intended to capture an idea about "length": namely, the length of a vector $\mathbf{v}$ should be the (inner) product of $\mathbf{v}$ with itself, and lengths are supposed to be nonnegative. Furthermore, the only vector of length zero should be the zero vector.

- Definition: A vector space $V$ together with an inner product $\langle \cdot, \cdot \rangle$ on $V$ is called an inner product space.

    ○ Any given vector space may have many different inner products.

    ○ When we say "Let $V$ be an inner product space", we intend this to mean that $V$ is equipped with a particular (fixed) inner product.

- The entire purpose of defining an inner product is to generalize the notion of the dot product to more general vector spaces, so we should check that the dot product on $\mathbb{R}^n$ is actually an inner product:

- Example: Show that the standard dot product on $\mathbb{R}^n$, defined as $(x_1, \ldots, x_n) \cdot (y_1, \ldots, y_n) = x_1 y_1 + \cdots + x_n y_n$ is an inner product.

    ○ [I1]-[I2]: We have previously verified the linearity and symmetry properties.

    ○ [I3]: If $\mathbf{v} = (x_1, \ldots, x_n)$ then $\mathbf{v} \cdot \mathbf{v} = x_1^2 + x_2^2 + \cdots + x_n^2$. Since each square is nonnegative, $\mathbf{v} \cdot \mathbf{v} \geq 0$, and $\mathbf{v} \cdot \mathbf{v} = 0$ only when all of the components of $\mathbf{v}$ are zero.

- There are other examples of inner products on $\mathbb{R}^n$ beyond the standard dot product.

- Example: Show that the pairing $\langle (x_1, y_1), (x_2, y_2) \rangle = 3x_1 x_2 + 2x_1 y_2 + 2x_2 y_1 + 4y_1 y_2$ on $\mathbb{R}^2$ is an inner product.

    ○ [I1]-[I2]: It is an easy algebraic computation to verify the linearity and symmetry properties.

    ○ [I3]: We have $\langle (x, y), (x, y) \rangle = 3x^2 + 4xy + 4y^2 = 2x^2 + (x + 2y)^2$, and since each square is nonnegative, the inner product is always nonnegative. Furthermore, it equals zero only when both squares are zero, and this clearly only occurs for $x = y = 0$.

- We can define another class of inner products on function spaces using integration:

- Example: Let $V$ be the vector space of continuous (real-valued) functions on the interval $[a, b]$. Show that $\langle f, g \rangle = \int_a^b f(x) g(x) \, dx$ is an inner product on $V$.

    ○ [I1]: We have $\langle f_1 + cf_2, g \rangle = \int_a^b [f_1(x) + cf_2(x)] \, g(x) \, dx = \int_a^b f_1(x) g(x) \, dx + c \int_a^b f_2(x) g(x) \, dx = \langle f_1, g \rangle + c \langle f_2, g \rangle$.

○ [I2]: Observe that $\langle g, f \rangle = \int_a^b g(x)f(x)\,dx = \int_a^b f(x)g(x)\,dx = \langle f, g \rangle$.

○ [I3]: Notice that $\langle f, f \rangle = \int_a^b f(x)^2\,dx$ is the integral of a nonnegative function, so it is always nonnegative. Furthermore (since $f$ is assumed to be continuous) the integral of $f^2$ cannot be zero unless $f$ is identically zero.

○ <u>Remark</u>: More generally, if $w(x)$ is any fixed positive ("weight") function that is continuous on $[a, b]$, $\langle f, g \rangle = \int_a^b f(x)g(x) \cdot w(x)\,dx$ is an inner product on $V$.

### 3.1.2 Inner Products on Complex Vector Spaces

- We would now like to extend the notion of an inner product to complex vector spaces.

  ○ A natural first guess would be to use the same definition as in a real vector space. However, this turns out not to be the right choice!

  ○ To explain why, suppose that we try to use the same definition of dot product to find the "length" of a vector of complex numbers, i.e., by computing $\mathbf{v} \cdot \mathbf{v}$.

  ○ We can see that the dot product of the vector $(1, 0, i)$ with itself is 0, but this vector certainly is not the zero vector. Even worse, the dot product of $(1, 0, 2i)$ with itself is $-3$, while the dot product of $(1, 0, 1 + 2i)$ with itself is $-2 + 4i$. None of these choices for "lengths" seem sensible.

  ○ Indeed, a more natural choice of "length" for the vector $(1, 0, i)$ is $\sqrt{2}$: the first component has absolute value 1, the second has absolute value 0, and the last has absolute value 1, for an overall "length" of $\sqrt{1^2 + 0^2 + 1^2}$, in analogy with the real vector $(1, 0, 1)$ which also has length $\sqrt{2}$. Similarly, $(1, 0, 2i)$ should really have length $\sqrt{1^2 + 0^2 + 2^2} = \sqrt{5}$.

  ○ One way to obtain a nonnegative function that seems to capture this idea of "length" for a complex vector is to include a conjugation: notice that for any complex vector $\mathbf{v}$, the dot product $\mathbf{v} \cdot \overline{\mathbf{v}}$ will always be a nonnegative real number.

  ○ Using the example above, we can compute that $(1, 0, i) \cdot \overline{(1, 0, i)} = 1^2 + 0^2 + 1^2 = 2$, so this "modified dot product" seems to give the square of the length of a complex vector, at least in this one case.

  ○ All of this suggests that the right analogy[1] of the "$\mathbb{R}^n$ dot product" for a pair of complex vectors $\mathbf{v}$, $\mathbf{w}$ is $\mathbf{v} \cdot \overline{\mathbf{w}}$ rather than $\mathbf{v} \cdot \mathbf{w}$.

  ○ We then lose the symmetry property of the real inner product, since now $\mathbf{v} \cdot \overline{\mathbf{w}}$ and $\mathbf{w} \cdot \overline{\mathbf{v}}$ are not equal in general. But notice we still have a relation between $\mathbf{v} \cdot \overline{\mathbf{w}}$ and $\mathbf{w} \cdot \overline{\mathbf{v}}$: namely, the latter is the complex conjugate of the former.

  ○ We can therefore obtain the right definition by changing the symmetry property $\langle \mathbf{w}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$ to "conjugate-symmetry": $\langle \mathbf{w}, \mathbf{v} \rangle = \overline{\langle \mathbf{v}, \mathbf{w} \rangle}$ .

- With this minor modification, we can extend the idea of an inner product to a complex vector space:

- <u>Definition</u>: If $V$ is a (real or) complex vector space, an <u>inner product</u> on $V$ is a pairing that assigns a scalar in $F$ to each ordered pair $(\mathbf{v}, \mathbf{w})$ of vectors in $V$. This pairing is denoted $\langle \mathbf{v}, \mathbf{w} \rangle$ and must satisfy the following properties:

  **[I1]** Linearity in the first argument: $\langle \mathbf{v}_1 + c\mathbf{v}_2, \mathbf{w} \rangle = \langle \mathbf{v}_1, \mathbf{w} \rangle + c\langle \mathbf{v}_2, \mathbf{w} \rangle$.

  **[I2]** Conjugate-symmetry: $\langle \mathbf{w}, \mathbf{v} \rangle = \overline{\langle \mathbf{v}, \mathbf{w} \rangle}$ (where the bar denotes the complex conjugate).

  **[I3]** Positive-definiteness: $\langle \mathbf{v}, \mathbf{v} \rangle \geq 0$ for all $\mathbf{v}$, and $\langle \mathbf{v}, \mathbf{v} \rangle = 0$ only when $\mathbf{v} = \mathbf{0}$.

  ○ This definition generalizes the one we gave for a real vector space earlier: if $V$ is a real vector space then $\overline{\langle \mathbf{w}, \mathbf{v} \rangle} = \langle \mathbf{w}, \mathbf{v} \rangle$ since $\langle \mathbf{w}, \mathbf{v} \rangle$ is always a real number.

---

[1] Another (perhaps more compelling) reason that symmetry is not the right property for complex vector spaces is that the only "symmetric positive-definite complex bilinear pairing", by which we mean a complex vector space $V$ with an inner product $\langle \cdot, \cdot \rangle$ satisfying the three axioms we listed for a real inner product space, is the trivial inner product on the zero space. Therefore, because we want to develop an interesting theory, we instead require conjugate-symmetry.

- ○ <u>Important Warning</u>: In other disciplines (particularly physics), inner products are often defined to be linear in the second argument, rather than the first. With this convention, the roles of the first and second component will be reversed (relative to our definition). This does not make any difference in the general theory, but can be extremely confusing since the definitions in mathematics and physics are otherwise identical, and the properties will have very similar statements.

- The chief reason for using this definition is that our modified dot product for complex vectors is an inner product:

- <u>Example</u>: For complex numbers $\mathbf{v} = (x_1, \ldots, x_n)$ and $\mathbf{w} = (y_1, \ldots, y_n)$, show that the map $\langle \mathbf{v}, \mathbf{w} \rangle = x_1\overline{y_1} + \cdots + x_n\overline{y_n}$ is an inner product on $\mathbb{C}^n$.

  - ○ [I1]: If $\mathbf{v}_1 = (x_1, \ldots, x_n)$, $\mathbf{v}_2 = (z_1, \ldots, z_n)$, and $\mathbf{w} = (y_1, \ldots, y_n)$, then

  $$\begin{aligned} \langle \mathbf{v}_1 + c\mathbf{v}_2, \mathbf{w} \rangle &= (x_1 + cz_1)\overline{y_1} + \cdots + (x_n + cz_n)\overline{y_n} \\ &= (x_1\overline{y_1} + \cdots + x_n\overline{y_n}) + c(z_1\overline{y_1} + \cdots + z_n\overline{y_n}) \\ &= \langle \mathbf{v}_1, \mathbf{w} \rangle + c\langle \mathbf{v}_2, \mathbf{w} \rangle. \end{aligned}$$

  - ○ [I2]: Observe that $\overline{\langle \mathbf{w}, \mathbf{v} \rangle} = \overline{y_1\overline{x_1} + \cdots + y_n\overline{x_n}} = \overline{y_1}x_1 + \cdots + \overline{y_n}x_n = \langle \mathbf{v}, \mathbf{w} \rangle$.

  - ○ [I3]: If $\mathbf{v} = (x_1, \ldots, x_n)$ then $\langle \mathbf{v}, \mathbf{v} \rangle = x_1\overline{x_1} + x_2\overline{x_2} + \cdots + x_n\overline{x_n} = |x_1|^2 + \cdots + |x_n|^2$. Each term is nonnegative, so $\mathbf{v} \cdot \mathbf{v} \geq 0$, and clearly $\mathbf{v} \cdot \mathbf{v} = 0$ only when all of the components of $\mathbf{v}$ are zero.

  - ○ This map is often called the "standard inner product" on $\mathbb{C}^n$, since it is fairly natural.

- Here is an example of a complex inner product on the space of $n \times n$ matrices:

- <u>Example</u>: Let $V = M_{n \times n}(\mathbb{C})$ be the vector space of complex $n \times n$ matrices. Show that $\langle A, B \rangle = \mathrm{tr}(AB^*)$ is an inner product on $V$, where $M^* = \overline{M^T}$ is the complex conjugate of the transpose of $M$ (often called the <u>conjugate transpose</u> or the <u>adjoint</u> of $M$).

  - ○ [I1]: We have $\langle A + cC, B \rangle = \mathrm{tr}[(A+cC)B^*] = \mathrm{tr}[AB^* + cCB^*] = \mathrm{tr}(AB^*) + c\,\mathrm{tr}(CB^*) = \langle A, B \rangle + c\,\langle C, B \rangle$, where we used the facts that $\mathrm{tr}(M + N) = \mathrm{tr}(M) + \mathrm{tr}(M)$ and $\mathrm{tr}(cM) = c\,\mathrm{tr}(M)$.

  - ○ [I2]: Observe that $\overline{\langle B, A \rangle} = \overline{\mathrm{tr}(BA^*)} = \mathrm{tr}(B^*A) = \mathrm{tr}(AB^*) = \langle A, B \rangle$, where we used the facts that $\mathrm{tr}(\overline{MN}) = \mathrm{tr}(M^*N^*)$ and that $\mathrm{tr}(MN) = \mathrm{tr}(NM)$, both of which are easy algebraic calculations.

  - ○ [I3]: We have $\langle A, A \rangle = \sum_{j=1}^{n}(AA^*)_{j,j} = \sum_{j=1}^{n}\sum_{k=1}^{n} A_{j,k}A_{k,j}^* = \sum_{j=1}^{n}\sum_{k=1}^{n} A_{j,k}\overline{A_{j,k}} = \sum_{j=1}^{n}\sum_{k=1}^{n} |A_{j,k}|^2 \geq 0$, and equality can only occur when each element of $A$ has absolute value zero (i.e., is zero).

  - ○ <u>Remark</u>: This inner product is often called the <u>Frobenius inner product</u>.

### 3.1.3 Properties of Inner Products, Norms

- Our fundamental goal in studying inner products is to extend the notion of length in $\mathbb{R}^n$ to a more general setting. Using the positive-definiteness property, we can define a notion of length in an inner product space.

- <u>Definition</u>: If $V$ is an inner product space, we define the <u>norm</u> (or <u>length</u>) of a vector $\mathbf{v}$ to be $||\mathbf{v}|| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$.

  - ○ When $V = \mathbb{R}^n$ with the standard dot product, the norm on $V$ reduces to the standard notion of "length" of a vector in $\mathbb{R}^n$.

- Here are a few basic properties of inner products and norms:

- <u>Proposition</u> (Properties of Norms): If $V$ is an inner product space with inner product $\langle \cdot, \cdot \rangle$, the following hold:

  1. For any vectors $\mathbf{v}$, $\mathbf{w}_1$, and $\mathbf{w}_2$, $\langle \mathbf{v}, \mathbf{w}_1 + \mathbf{w}_2 \rangle = \langle \mathbf{v}, \mathbf{w}_1 \rangle + \langle \mathbf{v}, \mathbf{w}_2 \rangle$.
     - ○ <u>Proof</u>: We have $\langle \mathbf{v}, \mathbf{w}_1 + \mathbf{w}_2 \rangle = \overline{\langle \mathbf{w}_1 + \mathbf{w}_2, \mathbf{v} \rangle} = \overline{\langle \mathbf{w}_1, \mathbf{v} \rangle + \langle \mathbf{w}_2, \mathbf{v} \rangle} = \langle \mathbf{v}, \mathbf{w}_1 \rangle + \langle \mathbf{v}, \mathbf{w}_2 \rangle$ by [I1] and [I2].
  2. For any vectors $\mathbf{v}$ and $\mathbf{w}$ and scalar $c$, $\langle \mathbf{v}, c\mathbf{w} \rangle = \bar{c}\,\langle \mathbf{v}, \mathbf{w} \rangle$.

- ○ Proof: We have $\langle \mathbf{v}, c\mathbf{w} \rangle = \overline{\langle c\mathbf{w}, \mathbf{v} \rangle} = \overline{c}\overline{\langle \mathbf{w}, \mathbf{v} \rangle} = \overline{c}\langle \mathbf{v}, \mathbf{w} \rangle$ by [I1] and [I2].

3. For any vector $\mathbf{v}$, $\langle \mathbf{v}, \mathbf{0} \rangle = 0 = \langle \mathbf{0}, \mathbf{v} \rangle$.
    - ○ Proof: Apply property (3) and [I2] with $c = 0$, using the fact that $0\mathbf{w} = \mathbf{0}$ for any $\mathbf{w}$.
4. For any vector $\mathbf{v}$, $||\mathbf{v}||$ is a nonnegative real number, and $||\mathbf{v}|| = 0$ if and only if $\mathbf{v} = \mathbf{0}$.
    - ○ Proof: Immediate from [I3].
5. For any vector $\mathbf{v}$ and scalar $c$, $||c\mathbf{v}|| = |c| \cdot ||\mathbf{v}||$.
    - ○ Proof: We have $||c \cdot \mathbf{v}|| = \sqrt{\langle c \cdot \mathbf{v}, c \cdot \mathbf{v} \rangle} = \sqrt{c\overline{c}\langle \mathbf{v}, \mathbf{v} \rangle} = |c| \cdot ||\mathbf{v}||$, using [I2] and property (2).

- In $\mathbb{R}^n$, there are a number of fundamental inequalities about lengths, which generalize quite pleasantly to general inner product spaces.

- The following result, in particular, is one of the most fundamental inequalities in all of mathematics:

- Theorem (Cauchy-Schwarz Inequality): For any $\mathbf{v}$ and $\mathbf{w}$ in an inner product space $V$, we have $|\langle \mathbf{v}, \mathbf{w} \rangle| \leq ||\mathbf{v}|| \, ||\mathbf{w}||$, with equality if and only if the set $\{\mathbf{v}, \mathbf{w}\}$ is linearly dependent.

    - ○ Proof: If $\mathbf{w} = \mathbf{0}$ then the result is trivial (since both sides are zero, and $\{\mathbf{v}, \mathbf{0}\}$ is always dependent), so now assume $\mathbf{w} \neq \mathbf{0}$.
    - ○ Let $t = \dfrac{\langle \mathbf{v}, \mathbf{w} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle}$. By properties of inner products and norms, we can write

$$
\begin{aligned}
||\mathbf{v} - t\mathbf{w}||^2 &= \langle \mathbf{v} - t\mathbf{w}, \mathbf{v} - t\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{v} \rangle - t\langle \mathbf{w}, \mathbf{v} \rangle - \overline{t}\langle \mathbf{v}, \mathbf{w} \rangle + t\overline{t}\langle \mathbf{w}, \mathbf{w} \rangle \\
&= \langle \mathbf{v}, \mathbf{v} \rangle - \frac{\overline{\langle \mathbf{v}, \mathbf{w} \rangle}\langle \mathbf{v}, \mathbf{w} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle} - \frac{|\langle \mathbf{v}, \mathbf{w} \rangle|^2}{\langle \mathbf{w}, \mathbf{w} \rangle} + \frac{|\langle \mathbf{v}, \mathbf{w} \rangle|^2}{\langle \mathbf{w}, \mathbf{w} \rangle} \\
&= \langle \mathbf{v}, \mathbf{v} \rangle - \frac{|\langle \mathbf{v}, \mathbf{w} \rangle|^2}{\langle \mathbf{w}, \mathbf{w} \rangle}.
\end{aligned}
$$

    - ○ Therefore, since $||\mathbf{v} - t\mathbf{w}||^2 \geq 0$ and $\langle \mathbf{w}, \mathbf{w} \rangle \geq 0$, clearing denominators and rearranging yields $|\langle \mathbf{v}, \mathbf{w} \rangle|^2 \leq \langle \mathbf{v}, \mathbf{v} \rangle \langle \mathbf{w}, \mathbf{w} \rangle$. Taking the square root yields the stated result.
    - ○ Furthermore, we will have equality if and only if $||\mathbf{v} - t\mathbf{w}||^2 = 0$, which is in turn equivalent to $\mathbf{v} - t\mathbf{w} = \mathbf{0}$; namely, when $\mathbf{v}$ is a multiple of $\mathbf{w}$. Since we also get equality if $\mathbf{w} = \mathbf{0}$, equality occurs precisely when the set $\{\mathbf{v}, \mathbf{w}\}$ is linearly dependent.
    - ○ Remark: As written, this proof is completely mysterious: why does making that particular choice for $t$ work? Here is some motivation: in the special case where $V$ is a real vector space, we can write $||\mathbf{v} - t\mathbf{w}|| = \langle \mathbf{v}, \mathbf{v} \rangle - 2t\langle \mathbf{v}, \mathbf{w} \rangle + t^2\langle \mathbf{w}, \mathbf{w} \rangle$, which is a quadratic function of $t$ that is always nonnegative.
    - ○ To decide whether a quadratic function is always nonnegative, we complete the square to see that
$$
\langle \mathbf{v}, \mathbf{v} \rangle - 2t\langle \mathbf{v}, \mathbf{w} \rangle + t^2\langle \mathbf{w}, \mathbf{w} \rangle = \langle \mathbf{w}, \mathbf{w} \rangle \left[ t - \frac{\langle \mathbf{v}, \mathbf{w} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle} \right]^2 + \left[ \langle \mathbf{v}, \mathbf{v} \rangle - \frac{\langle \mathbf{v}, \mathbf{w} \rangle^2}{\langle \mathbf{w}, \mathbf{w} \rangle} \right].
$$
    - ○ Thus, the minimum value of the quadratic function is $\langle \mathbf{v}, \mathbf{v} \rangle - \dfrac{\langle \mathbf{v}, \mathbf{w} \rangle^2}{\langle \mathbf{w}, \mathbf{w} \rangle}$, and it occurs when $t = \dfrac{\langle \mathbf{v}, \mathbf{w} \rangle}{\langle \mathbf{w}, \mathbf{w} \rangle}$.

- The Cauchy-Schwarz inequality has many applications (most of which are, naturally, proving other inequalities). Here are a few such applications:

- Theorem (Triangle Inequality): For any $\mathbf{v}$ and $\mathbf{w}$ in an inner product space $V$, we have $||\mathbf{v} + \mathbf{w}|| \leq ||\mathbf{v}|| + ||\mathbf{w}||$, with equality if and only if one vector is a positive-real scalar multiple of the other.

    - ○ Proof: Using the Cauchy-Schwarz inequality and the fact that $\text{Re}(z) \leq |z|$ for any $z \in \mathbb{C}$, we have

$$
\begin{aligned}
||\mathbf{v} + \mathbf{w}||^2 = \langle \mathbf{v} + \mathbf{w}, \mathbf{v} + \mathbf{w} \rangle &= \langle \mathbf{v}, \mathbf{v} \rangle + 2\text{Re}(\langle \mathbf{v}, \mathbf{w} \rangle) + \langle \mathbf{w}, \mathbf{w} \rangle \\
&\leq \langle \mathbf{v}, \mathbf{v} \rangle + 2|\langle \mathbf{v}, \mathbf{w} \rangle| + \langle \mathbf{w}, \mathbf{w} \rangle \\
&\leq \langle \mathbf{v}, \mathbf{v} \rangle + 2||\mathbf{v}|| \, ||\mathbf{w}|| + \langle \mathbf{w}, \mathbf{w} \rangle \\
&= ||\mathbf{v}||^2 + 2||\mathbf{v}|| \, ||\mathbf{w}|| + ||\mathbf{w}||^2.
\end{aligned}
$$

Taking the square root of both sides yields the desired result.

- ○ Equality will hold if and only if $\{\mathbf{v}, \mathbf{w}\}$ is linearly dependent (for equality in the Cauchy-Schwarz inequality) and $\langle \mathbf{v}, \mathbf{w} \rangle$ is a nonnegative real number. If either vector is zero, equality always holds. Otherwise, we must have $\mathbf{v} = c \cdot \mathbf{w}$ for some nonzero constant $c$: then $\langle \mathbf{v}, \mathbf{w} \rangle = c \langle \mathbf{w}, \mathbf{w} \rangle$ will be a nonnegative real number if and only if $c$ is a nonnegative real number.

- Example: Show that for any continuous function $f$ on $[0, 3]$, it is true that $\int_0^3 x f(x)\, dx \leq 3\sqrt{\int_0^3 f(x)^2\, dx}$.

  - ○ Simply apply the Cauchy-Schwarz inequality to $f$ and $g(x) = x$ in the inner product space of continuous functions on $[0, 3]$ with inner product $\langle f, g \rangle = \int_0^3 f(x) g(x)\, dx$.
  - ○ We obtain $|\langle f, g \rangle| \leq ||f|| \cdot ||g||$, or, explicitly, $\left| \int_0^3 x f(x)\, dx \right| \leq \sqrt{\int_0^3 f(x)^2\, dx} \cdot \sqrt{\int_0^3 x^2\, dx} = 3\sqrt{\int_0^3 f(x)^2\, dx}$.
  - ○ Since any real number is less than or equal to its absolute value, we immediately obtain the required inequality $\int_0^3 x f(x)\, dx \leq 3\sqrt{\int_0^3 f(x)^2\, dx}$.

- Example: Show that for any positive reals $a, b, c$, it is true that $\sqrt{\dfrac{a + 2b}{a + b + c}} + \sqrt{\dfrac{b + 2c}{a + b + c}} + \sqrt{\dfrac{c + 2a}{a + b + c}} \leq 3$.

  - ○ Let $\mathbf{v} = (\sqrt{a + 2b}, \sqrt{b + 2c}, \sqrt{c + 2a})$ and $\mathbf{w} = (1, 1, 1)$ in $\mathbb{R}^3$. By the Cauchy-Schwarz inequality, $\mathbf{v} \cdot \mathbf{w} \leq ||\mathbf{v}||\, ||\mathbf{w}||$.
  - ○ We compute $\mathbf{v} \cdot \mathbf{w} = \sqrt{a + 2b} + \sqrt{b + 2c} + \sqrt{c + 2a}$, along with $||\mathbf{v}||^2 = (a + 2b) + (b + 2c) + (c + 2a) = 3(a + b + c)$ and $||\mathbf{w}||^2 = 3$.
  - ○ Thus, we see $\sqrt{a + 2b} + \sqrt{a + 2c} + \sqrt{b + 2c} \leq \sqrt{3(a + b + c)} \cdot \sqrt{3}$, and upon dividing through by $\sqrt{a + b + c}$ we obtain the required inequality.

- Example (for those who like quantum mechanics): Prove the momentum-position formulation of Heisenberg's uncertainty principle: $\sigma_x \sigma_p \geq \overline{h}/2$. (In words: the product of uncertainties of position and momentum is greater than or equal to half of the reduced Planck constant.)

  - ○ It is a straightforward computation that, for two (complex-valued) observables $X$ and $Y$, the pairing $\langle X, Y \rangle = E[X\overline{Y}]$, the expected value of $X\overline{Y}$, is an inner product on the space of observables.
  - ○ Assume (for simplicity) that $x$ and $p$ both have expected value 0.
  - ○ We assume as given the commutation relation $xp - px = i\overline{h}$.
  - ○ By definition, $(\sigma_x)^2 = E[x\overline{x}] = \langle x, x \rangle$ and $(\sigma_p)^2 = E[p\overline{p}] = E[\overline{p}p] = \langle \overline{p}, \overline{p} \rangle$ are the variances of $x$ and $p$ respectively (in the statistical sense).
  - ○ By the Cauchy-Schwarz inequality, we can therefore write $\sigma_x^2 \sigma_p^2 = \langle x, x \rangle \langle \overline{p}, \overline{p} \rangle \geq |\langle x, \overline{p} \rangle|^2 = |E[xp]|^2$.
  - ○ We can write $xp = \dfrac{1}{2}(xp + px) + \dfrac{1}{2}(xp - px)$, where the first component is real and the second is imaginary, so taking expectations yields $E[xp] = \dfrac{1}{2}E[xp + px] + \dfrac{1}{2}E[xp - px]$, and therefore, $|E[xp]| \geq \dfrac{1}{2}|E[xp - px]| = \dfrac{1}{2}\left|i\overline{h}\right| = \dfrac{\overline{h}}{2}$.
  - ○ Combining with the inequality above yields $\sigma_x^2 \sigma_p^2 \geq \overline{h}^2/4$, and taking square roots yields $\sigma_x \sigma_p \geq \overline{h}/2$.

- Example (for those who like statistics): Prove the Cramér-Rao inequality: if $p(x; \theta)$ is a positive probability density function in $x$ on a finite sample space $S$ that is differentiable in $\theta$, and $g(x)$ is an unbiased estimator for the parameter $\theta$, then $E[(g(x) - \theta)^2] \geq 1/I(\theta)$ where $I : S \to \mathbb{R}$ is defined by $I(\theta) = E[(p_\theta(x; \theta)/p(x; \theta))^2]$. (In words: the variance of the estimator $g$ is greater than or equal to the reciprocal of the Fisher information of $p$.)

  - ○ By definition we have $\sum_{x \in S} p(x; \theta) = 1$ since $p$ is a probability density function on $S$. Also, since the expected value of a function $g : S \to \mathbb{R}$ is the sum $\sum_{x \in S} p(x; \theta)\, g(x)$, saying that $g$ is unbiased is the same as saying that $\sum_{x \in S} p(x; \theta)\, g(x) = \theta$ for all $\theta$.
  - ○ Differentiating these two equalities with respect to $\theta$ then yields $\sum_{x \in S} p_\theta(x; \theta) = 0$ and $\sum_{x \in S} p_\theta(x; \theta)\, g(x) = 1$, and then subtracting these two yields $\sum_{x \in S} p_\theta(x; \theta)\, [g(x) - \theta] = 1$.

- We may rewrite this last equality as $1 = \sum_{x \in S} \left( p(x;\theta)^{1/2} \left[ g(x) - \theta \right] \right) \left( p_\theta(x;\theta)/p(x;\theta)^{1/2} \right)$.
- Now apply the Cauchy-Schwarz inequality to both sides: if $\Sigma$ is the sum above, then we obtain $1^2 = \Sigma^2 \leq \left[ \sum_{x \in S} p(x;\theta)[g(x) - \theta]^2 \right] \cdot \left[ \sum_{x \in S} p_\theta(x;\theta)^2/p(x;\theta) \right]$.
- This immediately yields $\sum_{x \in S} p(x;\theta)[g(x) - \theta]^2 \geq 1/\sum_{x \in S} p(x;\theta)[p_\theta(x;\theta)/p(x;\theta)]^2$, which is equivalent to the claimed statement $E[(g(x) - \theta)^2] \geq 1/I(\theta)$.
- Remark 1: Note that the quantity $||g(x) - \theta||^2 = \sum_{x \in S} p(x;\theta)[g(x) - \theta]^2$ represents the variance of $g(x)$, whereas the lower bound $1/I(\theta)$ does not depend on $g$ (it only depends on the pdf $p$). Thus, the Cramér-Rao inequality gives a hard lower bound on how much variation $g$ must have in its estimate for $\theta$, if it correctly predicts the right value $\theta$ on average.
- Remark 2: The formulation above is stated only for a probability density function on a finite sample space, but the same proof works for arbitrary discrete or continuous random variables as long as the support of $p$ (i.e., the set on which $p$ is nonzero) does not depend on $\theta$.

## 3.2 Orthogonality

- Motivated by the Cauchy-Schwarz inequality, we can define a notion of angle between two nonzero vectors in a real inner product space:

- Definition: If $V$ is a real inner product space, we define the angle between two nonzero vectors $\mathbf{v}$ and $\mathbf{w}$ to be the real number $\theta$ in $[0, \pi]$ satisfying $\cos \theta = \dfrac{\langle \mathbf{v}, \mathbf{w} \rangle}{||\mathbf{v}|| \, ||\mathbf{w}||}$.

  - By the Cauchy-Schwarz inequality, the quotient on the right is a real number in the interval $[-1, 1]$, so there is exactly one such angle $\theta$.

- Example: Compute the angle between the vectors $\mathbf{v} = (3, -4, 5)$ and $\mathbf{w} = (1, 2, -2)$ under the standard dot product on $\mathbb{R}^3$.

  - We have $\mathbf{v} \cdot \mathbf{w} = -15$, $||\mathbf{v}|| = \sqrt{\mathbf{v} \cdot \mathbf{v}} = 5\sqrt{2}$, and $||\mathbf{w}|| = \sqrt{\mathbf{w} \cdot \mathbf{w}} = 3$.
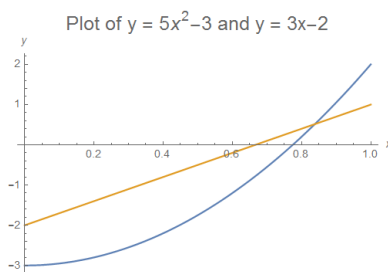  - Then the angle $\theta$ between the vectors satisfies $\cos(\theta) = \dfrac{-15}{15\sqrt{2}} = -\dfrac{1}{\sqrt{2}}$, so $\theta = \boxed{3\pi/4}$.

- Example: Compute the "angle" between $p = 5x^2 - 3$ and $q = 3x - 2$ in the inner product space of continuous functions on $[0, 1]$ with inner product $\langle f, g \rangle = \int_0^1 f(x)g(x)\,dx$.

  - We have $\langle p, q \rangle = \int_0^1 (5x^2 - 3)(3x - 2)\,dx = 23/12$, $||p|| = \sqrt{\int_0^1 (5x^2 - 3)^2\,dx} = 2$, and $||q|| = \sqrt{\int_0^1 (3x - 2)^2\,dx} = 1$.
  - Then the angle $\theta$ between the vectors satisfies $\cos(\theta) = \dfrac{23/12}{2} = \dfrac{23}{24}$, so $\theta = \boxed{\cos^{-1}(\dfrac{23}{24})}$.
  - The fact that this angle is so close to 0 suggests that these functions are nearly "parallel" in this inner product space. Indeed, the graphs of the two functions have very similar shapes:



Plot of y = 5x$^2$–3 and y = 3x–2

### 3.2.1 Orthogonality, Orthonormal Bases, and the Gram-Schmidt Procedure

- A particular case of interest is when the angle between two vectors is $\pi/2$ (i.e., they are perpendicular).

  - For nonzero vectors, by our discussion above, this is equivalent to saying that their inner product is 0.

- <u>Definition</u>: We say two vectors in an inner product space are <u>orthogonal</u> if their inner product is zero. We say a set $S$ of vectors is an <u>orthogonal set</u> if every pair of vectors in $S$ is orthogonal.

  - By our basic properties, the zero vector is orthogonal to every vector. Two nonzero vectors will be orthogonal if and only if the angle between them is $\pi/2$. (This generalizes the idea of two vectors being "perpendicular".)
  - <u>Example</u>: In $\mathbb{R}^3$ with the standard dot product, the vectors $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$ are orthogonal.
  - <u>Example</u>: In $\mathbb{R}^3$ with the standard dot product, the three vectors $(-1, 1, 2)$, $(2, 0, 1)$, and $(1, 5, -2)$ form an orthogonal set, since the dot product of each pair is zero.
  - The first orthogonal set above seems more natural than the second. One reason for this is that the vectors in the first set each have length 1, while the vectors in the second set have various different lengths ($\sqrt{6}$, $\sqrt{5}$, and $\sqrt{30}$ respectively).

- <u>Definition</u>: We say a set $S$ of vectors is an <u>orthonormal set</u> if every pair of vectors in $S$ is orthogonal, and every vector in $S$ has norm 1.

  - <u>Example</u>: In $\mathbb{R}^3$, $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ is an orthonormal set, but $\{(-1, 1, 2), (2, 0, 1), (1, 5, -2)\}$ is not.

- In both examples above, notice that the given orthogonal sets are also linearly independent. This feature is not an accident:

- <u>Proposition</u> (Orthogonality and Independence): In any inner product space, every orthogonal set of nonzero vectors is linearly independent.

  - <u>Proof</u>: Suppose we had a linear dependence $a_1 \mathbf{v}_1 + \cdots + a_n \mathbf{v}_n = \mathbf{0}$ for an orthogonal set of nonzero vectors $\{\mathbf{v}_1, \ldots, \mathbf{v}_n\}$.
  - Then, for any $j$, $0 = \langle \mathbf{0}, \mathbf{v}_j \rangle = \langle a_1 \mathbf{v}_1 + \cdots + a_n \mathbf{v}_n, \mathbf{v}_j \rangle = a_1 \langle \mathbf{v}_1, \mathbf{v}_j \rangle + \cdots + a_n \langle \mathbf{v}_n, \mathbf{v}_j \rangle = a_j \langle \mathbf{v}_j, \mathbf{v}_j \rangle$, since each of the inner products $\langle \mathbf{v}_i, \mathbf{v}_j \rangle$ for $i \neq j$ is equal to zero.
  - But now, since $\mathbf{v}_j$ is not the zero vector, $\langle \mathbf{v}_j, \mathbf{v}_j \rangle$ is positive, so it must be the case that $a_j = 0$. This holds for every $j$, so all the coefficients of the linear dependence are zero. Hence there can be no nontrivial linear dependence, so any orthogonal set is linearly independent.

- <u>Corollary</u>: If $V$ is an $n$-dimensional vector space and $S$ is an orthogonal set of $n$ nonzero vectors, then $S$ is a basis for $V$. (We refer to such a basis as an <u>orthogonal basis</u>.)

  - <u>Proof</u>: By the proposition above, $S$ is linearly independent, and by our earlier results, a linearly-independent set of $n$ vectors in an $n$-dimensional vector space is necessarily a basis.

- Given a basis of $V$, every vector in $V$ can be written as a unique linear combination of the basis vectors. However, actually computing the coefficients of the linear combination can be quite cumbersome.

  - If, however, we have an orthogonal basis for $V$, then we can compute the coefficients for the linear combination much more conveniently.

- <u>Theorem</u> (Orthogonal Decomposition): If $V$ is an $n$-dimensional vector space and $S = \{\mathbf{e}_1, \ldots, \mathbf{e}_n\}$ is an orthogonal basis, then for any $\mathbf{v}$ in $S$, we can write $\mathbf{v} = c_1 \mathbf{e}_1 + \cdots + c_n \mathbf{e}_n$, where $c_k = \dfrac{\langle \mathbf{v}, \mathbf{e}_k \rangle}{\langle \mathbf{e}_k, \mathbf{e}_k \rangle}$ for each $1 \leq k \leq n$. In particular, if $S$ is an orthonormal basis, then each $c_k = \langle \mathbf{v}, \mathbf{e}_k \rangle$.

  - <u>Proof</u>: Since $S$ is a basis, there do exist such coefficients $c_i$ and they are unique.

○ We then compute $\langle \mathbf{v}, \mathbf{e}_k \rangle = \langle c_1 \mathbf{e}_1 + \cdots + c_n \mathbf{e}_n, \mathbf{e}_k \rangle = c_1 \langle \mathbf{e}_1, \mathbf{e}_k \rangle + \cdots + c_n \langle \mathbf{e}_n, \mathbf{e}_k \rangle = c_k \langle \mathbf{e}_k, \mathbf{e}_k \rangle$ since each of the inner products $\langle \mathbf{e}_j, \mathbf{e}_k \rangle$ for $j \neq k$ is equal to zero.

○ Therefore, we must have $c_k = \dfrac{\langle \mathbf{v}, \mathbf{e}_k \rangle}{\langle \mathbf{e}_k, \mathbf{e}_k \rangle}$ for each $1 \leq k \leq n$.

○ If $S$ is an orthonormal basis, then $\langle \mathbf{e}_k, \mathbf{e}_k \rangle = 1$ for each $k$, so we get the simpler expression $c_k = \langle \mathbf{v}, \mathbf{e}_k \rangle$.

- <u>Example</u>: Write $\mathbf{v} = (7, 3, -4)$ as a linear combination of the basis $\{(-1, 1, 2), (2, 0, 1), (1, 5, -2)\}$ of $\mathbb{R}^3$.

    ○ We saw above that this set is an orthogonal basis, so let $\mathbf{e}_1 = (-1, 1, 2)$, $\mathbf{e}_2 = (2, 0, 1)$, and $\mathbf{e}_3 = (1, 5, -2)$.

    ○ We compute $\mathbf{v} \cdot \mathbf{e}_1 = -12$, $\mathbf{v} \cdot \mathbf{e}_2 = 10$, $\mathbf{v} \cdot \mathbf{e}_3 = 30$, $\mathbf{e}_1 \cdot \mathbf{e}_1 = 6$, $\mathbf{e}_2 \cdot \mathbf{e}_2 = 5$, and $\mathbf{e}_3 \cdot \mathbf{e}_3 = 30$.

    ○ Thus, per the theorem, $\mathbf{v} = c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 + c_3 \mathbf{e}_3$ where $c_1 = \dfrac{-12}{6} = -2$, $c_2 = \dfrac{10}{5} = 2$, and $c_3 = \dfrac{30}{30} = 1$.

    ○ Indeed, we can verify that $\boxed{(7, 3, -4) = -2(-1, 1, 2) + 2(2, 0, 1) + 1(1, 5, -2)}$.

- Given a basis, there exists a way to write any vector as a linear combination of the basis elements: the advantage of having an orthogonal basis is that we can easily *compute* the coefficients. We now give an algorithm for constructing an orthogonal basis for any finite-dimensional inner product space:

- <u>Theorem</u> (Gram-Schmidt Procedure): Let $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots\}$ be a basis of the inner product space $V$, and set $V_k = \mathrm{span}(\mathbf{v}_1, \dots, \mathbf{v}_k)$. Then there exists an orthogonal set of vectors $\{\mathbf{w}_1, \mathbf{w}_2, \dots\}$ such that, for each $k \geq 1$, $\mathrm{span}(\mathbf{w}_1, \dots, \mathbf{w}_k) = \mathrm{span}(V_k)$ and $\mathbf{w}_k$ is orthogonal to every vector in $V_{k-1}$. Furthermore, this sequence is unique up to multiplying the elements by nonzero scalars.

    ○ <u>Proof</u>: We construct the sequence $\{\mathbf{w}_1, \mathbf{w}_2, \dots\}$ recursively: we start with the simple choice $\mathbf{w}_1 = \mathbf{v}_1$.

    ○ Now suppose we have constructed $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{k-1}\}$, where $\mathrm{span}(\mathbf{w}_1, \dots, \mathbf{w}_{k-1}) = \mathrm{span}(V_{k-1})$.

    ○ Define the next vector as $\mathbf{w}_k = \mathbf{v}_k - a_1 \mathbf{w}_1 - a_2 \mathbf{w}_2 - \cdots - a_{k-1} \mathbf{w}_{k-1}$, where $a_j = \langle \mathbf{v}_k, \mathbf{w}_j \rangle / \langle \mathbf{w}_j, \mathbf{w}_j \rangle$.

    ○ From the construction, we can see that each of $\mathbf{w}_1, \dots, \mathbf{w}_k$ is a linear combination of $\mathbf{v}_1, \dots, \mathbf{v}_k$, and vice versa. Thus, by properties of span, $\mathrm{span}(\mathbf{w}_1, \dots, \mathbf{w}_k) = V_k$.

    ○ Then $\langle \mathbf{w}_k, \mathbf{w}_j \rangle = \langle \mathbf{v}_k - a_1 \mathbf{w}_1 - a_2 \mathbf{w}_2 - \cdots - a_{k-1} \mathbf{w}_{k-1}, \mathbf{w}_j \rangle = \langle \mathbf{v}_k, \mathbf{w}_j \rangle - a_1 \langle \mathbf{w}_1, \mathbf{w}_j \rangle - \cdots - a_{k-1} \langle \mathbf{w}_{k-1}, \mathbf{w}_j \rangle = \langle \mathbf{v}_k, \mathbf{w}_j \rangle - a_j \langle \mathbf{w}_j, \mathbf{w}_j \rangle = 0$ because all of the inner products $\langle \mathbf{w}_i, \mathbf{w}_j \rangle$ are zero except for $\langle \mathbf{w}_j, \mathbf{w}_j \rangle$.

    ○ Thus $\mathbf{w}_k$ is orthogonal to each of $\mathbf{w}_1, \dots, \mathbf{w}_{k-1}$ and hence also to all linear combinations of these vectors.

    ○ The uniqueness follows from the observation that (upon appropriate rescaling) we are essentially required to choose $\mathbf{w}_k = \mathbf{v}_k - a_1 \mathbf{w}_1 - a_2 \mathbf{w}_2 - \cdots - a_{k-1} \mathbf{w}_{k-1}$ for some scalars $a_1, \dots a_{k-1}$: orthogonality then forces the choice of the coefficients $a_j$ that we used above.

- <u>Corollary</u>: Every finite-dimensional inner product space has an orthonormal basis.

    ○ <u>Proof</u>: Choose any basis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ for $V$ and apply the Gram-Schmidt procedure: this yields an orthogonal basis $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ for $V$.

    ○ Now simply normalize each vector in $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$ by dividing by its norm: this preserves orthogonality, but rescales each vector to have norm 1, thus yielding an orthonormal basis for $V$.

- The proof of the Gram-Schmidt procedure may seem involved, but applying it in practice is fairly straightforward (if somewhat cumbersome computationally).

    ○ We remark here that, although our algorithm above gives an orthogonal basis, it is also possible to perform the normalization at each step during the procedure, to construct an orthonormal basis one vector at a time.

    ○ When performing computations by hand, it is generally disadvantageous to normalize at each step, because the norm of a vector will often involve square roots (which will then be carried into subsequent steps of the computation).

    ○ When using a computer (with approximate arithmetic), however, normalizing at each step can avoid certain numerical instability issues. The particular description of the algorithm we have discussed turns out not to be especially numerically stable, but it is possible to modify the algorithm to avoid magnifying the error as substantially when iterating the procedure.

- **Example:** For $V = \mathbb{R}^3$ with the standard inner product, apply the Gram-Schmidt procedure to the vectors $\mathbf{v}_1 = (2, 1, 2)$, $\mathbf{v}_2 = (5, 4, 2)$, $\mathbf{v}_3 = (-1, 2, 1)$. Use the result to find an orthonormal basis for $\mathbb{R}^3$.

  - We start with $\mathbf{w}_1 = \mathbf{v}_1 = \boxed{(2, 1, 2)}$.

  - Next, $\mathbf{w}_2 = \mathbf{v}_2 - a_1 \mathbf{w}_1$, where $a_1 = \dfrac{\mathbf{v}_2 \cdot \mathbf{w}_1}{\mathbf{w}_1 \cdot \mathbf{w}_1} = \dfrac{(5, 4, 2) \cdot (2, 1, 2)}{(2, 1, 2) \cdot (2, 1, 2)} = \dfrac{18}{9} = 2$. Thus, $\mathbf{w}_2 = \boxed{(1, 2, -2)}$.

  - Finally, $\mathbf{w}_3 = \mathbf{v}_3 - b_1 \mathbf{w}_1 - b_2 \mathbf{w}_2$ where $b_1 = \dfrac{\mathbf{v}_3 \cdot \mathbf{w}_1}{\mathbf{w}_1 \cdot \mathbf{w}_1} = \dfrac{(-1, 2, 1) \cdot (2, 1, 2)}{(2, 1, 2) \cdot (2, 1, 2)} = \dfrac{2}{9}$, and $b_2 = \dfrac{\mathbf{v}_3 \cdot \mathbf{w}_2}{\mathbf{w}_2 \cdot \mathbf{w}_2} = \dfrac{(-1, 2, 1) \cdot (1, 2, -2)}{(1, 2, -2) \cdot (1, 2, -2)} = \dfrac{1}{9}$. Thus, $\mathbf{w}_3 = \boxed{(-\dfrac{14}{9}, \dfrac{14}{9}, \dfrac{7}{9})}$.

  - For the orthonormal basis, we simply divide each vector by its length.

  - We get $\dfrac{\mathbf{w}_1}{||\mathbf{w}_1||} = (\dfrac{2}{3}, \dfrac{1}{3}, \dfrac{2}{3})$, $\dfrac{\mathbf{w}_2}{||\mathbf{w}_2||} = (\dfrac{1}{3}, \dfrac{2}{3}, -\dfrac{2}{3})$, and $\dfrac{\mathbf{w}_3}{||\mathbf{w}_3||} = (-\dfrac{2}{3}, \dfrac{2}{3}, \dfrac{1}{3})$.

- **Example:** For $V = \mathbb{R}[x]$ with inner product $\langle f, g \rangle = \int_0^1 f(x)g(x)\,dx$, apply the Gram-Schmidt procedure to the polynomials $p_1 = 1$, $p_2 = x$, $p_3 = x^2$, $p_4 = x^3$.

  - We start with $\mathbf{w}_1 = p_1 = \boxed{1}$.

  - Next, $\mathbf{w}_2 = p_2 - a_1 \mathbf{w}_1$, where $a_1 = \dfrac{\langle p_2, \mathbf{w}_1 \rangle}{\langle \mathbf{w}_1, \mathbf{w}_1 \rangle} = \dfrac{\int_0^1 x\,dx}{\int_0^1 1\,dx} = \dfrac{1}{2}$. Thus, $\mathbf{w}_2 = \boxed{x - \dfrac{1}{2}}$.

  - Then, $\mathbf{w}_3 = p_3 - b_1 \mathbf{w}_1 - b_2 \mathbf{w}_2$ where $b_1 = \dfrac{\langle p_3, \mathbf{w}_1 \rangle}{\langle \mathbf{w}_1, \mathbf{w}_1 \rangle} = \dfrac{\int_0^1 x^2\,dx}{\int_0^1 1\,dx} = \dfrac{1}{3}$, and $b_2 = \dfrac{\langle p_3, \mathbf{w}_2 \rangle}{\langle \mathbf{w}_2, \mathbf{w}_2 \rangle} = \dfrac{\int_0^1 x^2(x - 1/2)\,dx}{\int_0^1 (x - 1/2)^2\,dx} = \dfrac{1/12}{1/12} = 1$. Thus, $\mathbf{w}_3 = \boxed{x^2 - x + \dfrac{1}{6}}$.

  - Finally, $\mathbf{w}_4 = p_4 - c_1 \mathbf{w}_1 - c_2 \mathbf{w}_2 - c_3 \mathbf{w}_3$ where $b_1 = \dfrac{\langle p_4, \mathbf{w}_1 \rangle}{\langle \mathbf{w}_1, \mathbf{w}_1 \rangle} = \dfrac{\int_0^1 x^3\,dx}{\int_0^1 1\,dx} = \dfrac{1}{4}$, $b_2 = \dfrac{\langle p_4, \mathbf{w}_2 \rangle}{\langle \mathbf{w}_2, \mathbf{w}_2 \rangle} = \dfrac{\int_0^1 x^3(x - 1/2)\,dx}{\int_0^1 (x - 1/2)^2\,dx} = \dfrac{3/40}{1/12} = \dfrac{9}{10}$, and $b_3 = \dfrac{\langle p_4, \mathbf{w}_3 \rangle}{\langle \mathbf{w}_3, \mathbf{w}_3 \rangle} = \dfrac{\int_0^1 x^3(x^2 - x + 1/6)\,dx}{\int_0^1 (x^2 - x + 1/6)^2\,dx} = \dfrac{1/120}{1/180} = \dfrac{3}{2}$. Thus, $\mathbf{w}_4 = \boxed{x^3 - \dfrac{3}{2}x^2 + \dfrac{3}{5}x - \dfrac{1}{20}}$.

- We will mention that, although the Gram-Schmidt procedure allows us to construct an orthogonal basis for an arbitrary finite-dimensional vector space, there exist infinite-dimensional vector spaces that have no orthogonal basis.

  - The precise details are somewhat involved, but in fact, the space $\ell^2(\mathbb{R})$ of infinite real sequences $(a_1, a_2, \dots)$ such that $a_1^2 + a_2^2 + \cdots$ is finite, with inner product $\langle (a_1, a_2, \dots), (b_1, b_2, \dots) \rangle = a_1 b_1 + a_2 b_2 + \cdots$, has no orthogonal basis.

  - Notice, for example, that the set $\{e_1, e_2, \dots\}$ where $e_i$ has a 1 in the $i$th coordinate and 0s elsewhere is an orthonormal set but is not a basis. Nevertheless, this orthonormal set is maximal, in the sense that the only other element of the space orthogonal to every vector in the set is the zero vector: it cannot be extended to any larger orthonormal set.

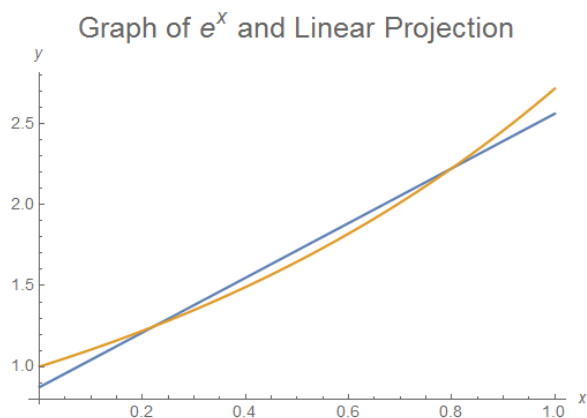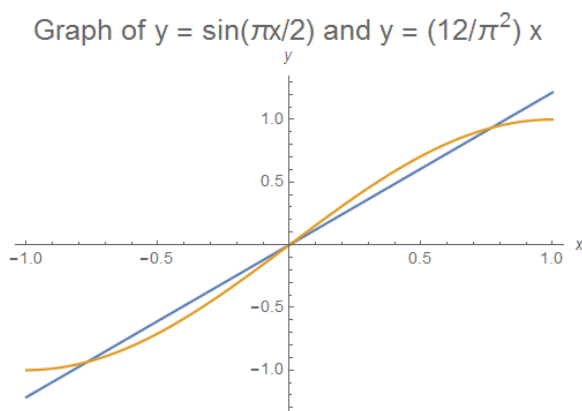### 3.2.2 Orthogonal Complements and Orthogonal Projection

- If $V$ is an inner product space, $W$ is a subspace, and $\mathbf{v}$ is some vector in $V$, we would like to study the problem of finding a "best approximation" of $\mathbf{v}$ in $W$.

  - For two vectors $\mathbf{v}$ and $\mathbf{w}$, the distance between $\mathbf{v}$ and $\mathbf{w}$ is $||\mathbf{v} - \mathbf{w}||$, so what we are seeking is a vector $\mathbf{w}$ in $S$ that minimizes the quantity $||\mathbf{v} - \mathbf{w}||$.

○ As a particular example, suppose we are given a point $P$ in $\mathbb{R}^2$ and wish to find the minimal distance from $P$ to a particular line in $\mathbb{R}^2$. Geometrically, the minimal distance is achieved by the segment $PQ$, where $Q$ is chosen so that $PQ$ is perpendicular to the line.

○ In a similar way, the minimal distance between a point in $\mathbb{R}^3$ and a given plane will also be minimized by finding the segment perpendicular to the plane.

○ Both of these problems suggest that the solution to this optimization problem will involve some notion of "perpendicularity" to the subspace $W$.

- <u>Definition</u>: Let $V$ be an inner product space. If $S$ is a nonempty subset of $V$, we say a vector $\mathbf{v}$ in $V$ is <u>orthogonal to $S$</u> if it is orthogonal to every vector in $S$. The set of all vectors orthogonal to $S$ is denoted $S^\perp$ ("$S$-perpendicular", or often "$S$-perp" for short).

  ○ We will typically be interested in the case where $S$ is a subspace of $V$. It is easy to see via the subspace criterion that $S^\perp$ is always a subspace of $V$, even if $S$ itself is not.

  ○ <u>Example</u>: In $\mathbb{R}^3$, if $W$ is the $xy$-plane consisting of all vectors of the form $(x, y, 0)$, then $W^\perp$ is the $z$-axis, consisting of the vectors of the form $(0, 0, z)$.

  ○ <u>Example</u>: In $\mathbb{R}^3$, if $W$ is the $x$-axis consisting of all vectors of the form $(x, 0, 0)$, then $W^\perp$ is the $yz$-plane, consisting of the vectors of the form $(0, y, z)$.

  ○ <u>Example</u>: In any inner product space $V$, $V^\perp = \{\mathbf{0}\}$ and $\{\mathbf{0}\}^\perp = V$.

- When $V$ is finite-dimensional, we can use the Gram-Schmidt process to compute an explicit basis of $W^\perp$:

- <u>Theorem</u> (Basis for Orthogonal Complement): Suppose $W$ is a subspace of the finite-dimensional inner product space $V$, and that $S = \{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$ is an orthonormal basis for $W$. If $\{\mathbf{e}_1, \ldots, \mathbf{e}_k, \mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$ is any extension of $S$ to an orthonormal basis for $V$, the set $\{\mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$ is an orthonormal basis for $W^\perp$. In particular, $\dim(V) = \dim(W) + \dim(W^\perp)$.

  ○ <u>Remark</u>: It is always possible to extend the orthonormal basis $S = \{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$ to an orthonormal basis for $V$: simply extend the linearly independent set $S$ to a basis $\{\mathbf{e}_1, \ldots, \mathbf{e}_k, \mathbf{x}_{k+1}, \ldots, \mathbf{x}_n\}$ of $V$, and then apply Gram-Schmidt to obtain an orthonormal basis $\{\mathbf{e}_1, \ldots, \mathbf{e}_k, \mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$.

  ○ <u>Proof</u>: For the first statement, the set $\{\mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$ is orthonormal and hence linearly independent. Since each vector is orthogonal to every vector in $S$, each of $\mathbf{e}_{k+1}, \ldots, \mathbf{e}_n$ is in $W^\perp$, and so it remains to show that $\{\mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$ spans $W^\perp$.

  ○ So let $\mathbf{v}$ be any vector in $W^\perp$. Since $\{\mathbf{e}_1, \ldots, \mathbf{e}_k, \mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$ is an orthonormal basis of $V$, by the orthogonal decomposition we know that $\mathbf{v} = \langle \mathbf{v}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \cdots + \langle \mathbf{v}, \mathbf{e}_k \rangle \mathbf{e}_k + \langle \mathbf{v}, \mathbf{e}_{k+1} \rangle \mathbf{e}_{k+1} + \cdots + \langle \mathbf{v}, \mathbf{e}_n \rangle \mathbf{e}_n$.

  ○ But since $\mathbf{v}$ is in $W^\perp$, $\langle \mathbf{v}, \mathbf{e}_1 \rangle = \cdots = \langle \mathbf{v}, \mathbf{e}_k \rangle = 0$: thus, $\mathbf{v} = \langle \mathbf{v}, \mathbf{e}_{k+1} \rangle \mathbf{e}_{k+1} + \cdots + \langle \mathbf{v}, \mathbf{e}_n \rangle \mathbf{e}_n$, and therefore $\mathbf{v}$ is contained in the span of $\{\mathbf{e}_{k+1}, \ldots, \mathbf{e}_n\}$, as required.

  ○ The statement about dimensions follows immediately from our explicit construction of the basis of $V$ as a union of the basis for $W$ and the basis for $W^\perp$.

- <u>Example</u>: If $W = \text{span}[\frac{1}{3}(1, 2, -2), \frac{1}{3}(-2, 2, 1)]$ in $\mathbb{R}^3$ with the standard dot product, find a basis for $W^\perp$.

  ○ Notice that the vectors $\mathbf{e}_1 = \frac{1}{3}(1, 2, -2)$ and $\mathbf{e}_2 = \frac{1}{3}(-2, 2, 1)$ form an orthonormal basis for $W$.

  ○ It is straightforward to verify that if $\mathbf{v}_3 = (1, 0, 0)$, then $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{v}_3\}$ is a linearly independent set and therefore a basis for $\mathbb{R}^3$.

  ○ Applying Gram-Schmidt to the set $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{v}_3\}$ yields $\mathbf{w}_1 = \mathbf{e}_1$, $\mathbf{w}_2 = \mathbf{e}_2$, and $\mathbf{w}_3 = \mathbf{v}_3 - \langle \mathbf{v}_3, \mathbf{w}_1 \rangle \mathbf{w}_1 - \langle \mathbf{v}_3, \mathbf{w}_2 \rangle \mathbf{w}_2 = \frac{1}{9}(4, 2, 4)$.

  ○ Normalizing $\mathbf{w}_3$ produces the orthonormal basis $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$ for $V$, with $\mathbf{e}_3 = \frac{1}{3}(2, 1, 2)$.

  ○ By the theorem above, we conclude that $\{\mathbf{e}_3\} = \boxed{\{\frac{1}{3}(2, 1, 2)\}}$ is an orthonormal basis of $W^\perp$.

○ Alternatively, we could have computed a basis for $W^\perp$ by observing that $\dim(W^\perp) = \dim(V) - \dim(W) = 1$, and then simply found one nonzero vector orthogonal to both $\frac{1}{3}(1, 2, -2)$ and $\frac{1}{3}(-2, 2, 1)$. (For this, we could have either solved the system of equations explicitly, or computed the cross product of the two given vectors.)

- We can give a simpler (although ultimately equivalent) method for finding a basis for $W^\perp$ using matrices:

- <u>Theorem</u> (Orthogonal Complements and Matrices): If $A$ is an $m \times n$ matrix, then the rowspace of $A$ and the nullspace of $A$ are orthogonal complements of one another in $\mathbb{R}^n$, with respect to the standard dot product.

  ○ <u>Proof</u>: Let $A$ be an $m \times n$ matrix, so that the rowspace and nullspace are both subspaces of $\mathbb{R}^n$.
  ○ By definition, any vector in rowspace($A$) is orthogonal to any vector in nullspace($A$), so rowspace($A$) $\subseteq$ nullspace($A$)$^\perp$ and nullspace($A$) $\subseteq$ rowspace($A$)$^\perp$.
  ○ Furthermore, since $\dim(\text{rowspace}(A)) + \dim(\text{nullspace}(A)) = n$ from our results on the respective dimensions of these spaces, we see that $\dim(\text{rowspace}(A)) = \dim(\text{nullspace}(A)^\perp)$ and $\dim(\text{nullspace}(A)) = \dim(\text{rowspace}(A)^\perp)$.
  ○ Since all these spaces are finite-dimensional, we must therefore have equality everywhere, as claimed.

- From the theorem above, when $W$ is a subspace of $\mathbb{R}^n$ with respect to the standard dot product, we can easily compute a basis for $W^\perp$ by computing the nullspace of the matrix whose rows are a spanning set for $W$.

  ○ Although this method is much faster, it will not produce an orthonormal basis of $W$. It can also be adapted for subspaces of an arbitrary finite-dimensional inner product space, but this requires having an orthonormal basis for the space computed ahead of time.

- <u>Example</u>: If $W = \text{span}[(1, 1, -1, 1), (1, 2, 0, -2)]$ in $\mathbb{R}^4$ with the standard dot product, find a basis for $W^\perp$.

  ○ We row-reduce the matrix whose rows are the given basis for $W$:

$$\begin{bmatrix} 1 & 1 & -1 & 1 \\ 1 & 2 & 0 & -2 \end{bmatrix} \xrightarrow{R_2 - R_1} \begin{bmatrix} 1 & 1 & -1 & 1 \\ 0 & 1 & 1 & -3 \end{bmatrix} \xrightarrow{R_1 - R_2} \begin{bmatrix} 1 & 0 & -2 & 4 \\ 0 & 1 & 1 & -3 \end{bmatrix}.$$

  ○ From the reduced row-echelon form, we see that $\boxed{\{(-4, 3, 0, 1), (2, -1, 1, 0)\}}$ is a basis for the nullspace and hence of $W^\perp$.

- As we might expect from geometric intuition, if $W$ is a subspace of the (finite-dimensional) inner product space $V$, we can decompose any vector uniquely as the sum of a component in $W$ with a component in $W^\perp$:

- <u>Theorem</u> (Orthogonal Components): Let $V$ be an inner product space and $W$ be a finite-dimensional subspace. Then every vector $\mathbf{v} \in V$ can be uniquely written in the form $\mathbf{v} = \mathbf{w} + \mathbf{w}^\perp$ for some $\mathbf{w} \in W$ and $\mathbf{w}^\perp \in W^\perp$, and furthermore, we have the Pythagorean relation $||\mathbf{v}||^2 = ||\mathbf{w}||^2 + ||\mathbf{w}^\perp||^2$.

  ○ <u>Proof</u>: First, we show that such a decomposition exists. Since $W$ is finite-dimensional, it has some orthonormal basis $\{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$.
  ○ Now set $\mathbf{w} = \langle \mathbf{v}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \langle \mathbf{v}, \mathbf{e}_2 \rangle \mathbf{e}_2 + \cdots + \langle \mathbf{v}, \mathbf{e}_k \rangle \mathbf{e}_k$, and then $\mathbf{w}^\perp = \mathbf{v} - \mathbf{w}$.
  ○ Clearly $\mathbf{w} \in W$ and $\mathbf{v} = \mathbf{w} + \mathbf{w}^\perp$, so we need only check that $\mathbf{w}^\perp \in W^\perp$.
  ○ To see this, first observe that $\langle \mathbf{w}, \mathbf{e}_i \rangle = \langle \mathbf{v}, \mathbf{e}_i \rangle$ since $\{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$ is an orthonormal basis. Then, we see that $\langle \mathbf{w}^\perp, \mathbf{e}_i \rangle = \langle \mathbf{v} - \mathbf{w}, \mathbf{e}_i \rangle = \langle \mathbf{v}, \mathbf{e}_i \rangle - \langle \mathbf{w}, \mathbf{e}_i \rangle = 0$. Thus, $\mathbf{w}^\perp$ is orthogonal to each vector in the orthonormal basis of $W$, so it is in $W^\perp$.
  ○ For the uniqueness, suppose we had two decompositions $\mathbf{v} = \mathbf{w}_1 + \mathbf{w}_1^\perp$ and $\mathbf{v} = \mathbf{w}_2 + \mathbf{w}_2^\perp$.
  ○ By subtracting and rearranging, we see that $\mathbf{w}_1 - \mathbf{w}_2 = \mathbf{w}_2^\perp - \mathbf{w}_1^\perp$. Denoting this common vector by $\mathbf{x}$, we see that $\mathbf{x}$ is in both $W$ and $W^\perp$: thus, $\mathbf{x}$ is orthogonal to itself, but the only such vector is the zero vector. Thus, $\mathbf{w}_1 = \mathbf{w}_2$ and $\mathbf{w}_1^\perp = \mathbf{w}_2^\perp$, so the decomposition is unique.
  ○ For the last statement, since $\langle \mathbf{w}, \mathbf{w}^\perp \rangle = 0$, we have $||\mathbf{v}||^2 = \langle \mathbf{v}, \mathbf{v} \rangle = \langle \mathbf{w} + \mathbf{w}^\perp, \mathbf{w} + \mathbf{w}^\perp \rangle = \langle \mathbf{w}, \mathbf{w} \rangle + \langle \mathbf{w}^\perp, \mathbf{w}^\perp \rangle = ||\mathbf{w}||^2 + ||\mathbf{w}^\perp||^2$, as claimed.

- <u>Definition</u>: If $V$ is an inner product space and $W$ is a finite-dimensional subspace with orthonormal basis $\{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$, the <u>orthogonal projection of $\mathbf{v}$ into $W$</u> is the vector $\operatorname{proj}_W(\mathbf{v}) = \langle \mathbf{v}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \langle \mathbf{v}, \mathbf{e}_2 \rangle \mathbf{e}_2 + \cdots + \langle \mathbf{v}, \mathbf{e}_k \rangle \mathbf{e}_k$.

    - If instead we only have an orthogonal basis $\{\mathbf{u}_1, \ldots, \mathbf{u}_k\}$ of $W$, the corresponding expression is instead $\operatorname{proj}_W(\mathbf{v}) = \dfrac{\langle \mathbf{v}, \mathbf{u}_1 \rangle}{\langle \mathbf{u}_1, \mathbf{u}_1 \rangle} \mathbf{u}_1 + \dfrac{\langle \mathbf{v}, \mathbf{u}_2 \rangle}{\langle \mathbf{u}_2, \mathbf{u}_2 \rangle} \mathbf{u}_2 + \cdots + \dfrac{\langle \mathbf{v}, \mathbf{u}_2 \rangle}{\langle \mathbf{u}_k, \mathbf{u}_k \rangle} \mathbf{u}_k$.

- <u>Example</u>: For $W = \operatorname{span}[(1,0,0), \frac{1}{5}(0,3,4)]$ in $\mathbb{R}^3$ under the standard dot product, compute the orthogonal projection of $\mathbf{v} = (1,2,1)$ into $W$, and verify the relation $||\mathbf{v}||^2 = ||\mathbf{w}||^2 + ||\mathbf{w}^\perp||^2$.

    - Notice that the vectors $\mathbf{e}_1 = (1,0,0)$ and $\mathbf{e}_2 = \frac{1}{5}(0,3,4)$ form an orthonormal basis for $W$.

    - Thus, the orthogonal projection is $\mathbf{w} = \operatorname{proj}_W(\mathbf{v}) = \langle \mathbf{v}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \langle \mathbf{v}, \mathbf{e}_2 \rangle \mathbf{e}_2 = 1\,(1,0,0) + 2\,(0,3/5,4/5) = \boxed{(1, 6/5, 8/5)}$.

    - We see that $\mathbf{w}^\perp = \mathbf{v} - \mathbf{w} = (0, 4/5, -3/5)$ is orthogonal to both $(1,0,0)$ and $(0,3/5,4/5)$, so it is indeed in $W^\perp$. Furthermore, $||\mathbf{v}||^2 = 6$, while $||\mathbf{w}||^2 = 5$ and $||\mathbf{w}^\perp||^2 = 1$, so indeed $||\mathbf{v}||^2 = ||\mathbf{w}||^2 + ||\mathbf{w}^\perp||^2$.

- The orthogonal projection gives the answer to the approximation problem we posed earlier:

- <u>Corollary</u> (Best Approximations): If $W$ is a finite-dimensional subspace of the inner product space $V$, then for any vector $\mathbf{v}$ in $V$, the orthogonal projection of $\mathbf{v}$ into $W$ is closer to $\mathbf{v}$ than any other vector in $W$. Explicitly, if $\mathbf{w}$ is the projection, then for any $\mathbf{w}' \in W$, we have $||\mathbf{v} - \mathbf{w}|| \leq ||\mathbf{v} - \mathbf{w}'||$ with equality if and only if $\mathbf{w} = \mathbf{w}'$.

    - <u>Proof</u>: By the theorem on orthogonal complements, we can write $\mathbf{v} = \mathbf{w} + \mathbf{w}^\perp$ where $\mathbf{w} \in W$ and $\mathbf{w}^\perp \in W^\perp$. Now, for any other vector $\mathbf{w}' \in W$, we can write $\mathbf{v} - \mathbf{w}' = (\mathbf{v} - \mathbf{w}) + (\mathbf{w} - \mathbf{w}')$, and observe that $\mathbf{v} - \mathbf{w} = \mathbf{w}^\perp$ is in $W^\perp$, and $\mathbf{w} - \mathbf{w}'$ is in $W$ (since both $\mathbf{w}$ and $\mathbf{w}'$ are, and $W$ is a subspace).

    - Thus, $\mathbf{v} - \mathbf{w}' = (\mathbf{v} - \mathbf{w}) + (\mathbf{w} - \mathbf{w}')$ is a decomposition of $\mathbf{v} - \mathbf{w}'$ into orthogonal vectors. Taking norms, we see that $||\mathbf{v} - \mathbf{w}'||^2 = ||\mathbf{v} - \mathbf{w}||^2 + ||\mathbf{w} - \mathbf{w}'||^2$.

    - Then, if $\mathbf{w}' \neq \mathbf{w}$, since the norm of $||\mathbf{w} - \mathbf{w}'||$ is positive, we conclude that $||\mathbf{v} - \mathbf{w}|| < ||\mathbf{v} - \mathbf{w}'||$.

- <u>Example</u>: Find the best approximation to $\mathbf{v} = (3, -3, 3)$ lying in the subspace $W = \operatorname{span}[\frac{1}{3}(1,2,-2), \frac{1}{3}(-2,2,1)]$, where distance is measured under the standard dot product.

    - Notice that the vectors $\mathbf{e}_1 = \frac{1}{3}(1,2,-2)$ and $\mathbf{e}_2 = \frac{1}{3}(-2,2,1)$ form an orthonormal basis for $W$.

    - Thus, the desired vector, the orthogonal projection, is $\operatorname{proj}_W(\mathbf{v}) = \langle \mathbf{v}, \mathbf{e}_1 \rangle \mathbf{e}_1 + \langle \mathbf{v}, \mathbf{e}_2 \rangle \mathbf{e}_2 = -3\mathbf{e}_1 - 3\mathbf{e}_2 = \boxed{(1, -4, 1)}$.

- <u>Example</u>: Find the best approximation to the function $f(x) = \sin(\pi x/2)$ that lies in the subspace $P_2(\mathbb{R})$, under the inner product $\langle f, g \rangle = \int_{-1}^{1} f(x) g(x)\, dx$.

    - First, by applying Gram-Schmidt to the basis $\{1, x, x^2\}$, we can generate an orthogonal basis of $P_2(\mathbb{R})$ under this inner product: the result (after rescaling to eliminate denominators) is $\{1, x, 3x^2 - 1\}$.

    - Now, with $p_1 = 1$, $p_2 = x$, $p_3 = 3x^2 - 1$ we can compute $\langle f, p_1 \rangle = 0$, $\langle f, p_2 \rangle = 8/\pi^2$, $\langle f, p_3 \rangle = 0$, and also $\langle p_1, p_1 \rangle = 2$, $\langle p_2, p_2 \rangle = 2/3$, and $\langle p_3, p_3 \rangle = 8/5$.

    - Thus, the desired orthogonal projection is $\operatorname{proj}_{P_1(\mathbb{R})}(f) = \dfrac{\langle f, p_1 \rangle}{\langle p_1, p_1 \rangle} p_1 + \dfrac{\langle f, p_2 \rangle}{\langle p_2, p_2 \rangle} p_2 + \dfrac{\langle f, p_3 \rangle}{\langle p_3, p_3 \rangle} p_3 = \boxed{\dfrac{12}{\pi^2} x}$.

    - We can see from a plot of the orthogonal projection polynomial and $f$ on $[-1, 1]$ that this line is a reasonably accurate approximation of $\sin(x)$ on the interval $[-1, 1]$:

Graph of y = sin(πx/2) and y = (12/π²) x

Graph of eˣ and Linear Projection

- <u>Example</u>: Find the linear polynomial $p(x)$ that minimizes the expression $\int_0^1 (p(x) - e^x)^2 \, dx$.

  - ○ Observe that the minimization problem is asking us to find the orthogonal projection of $e^x$ into $P_1(\mathbb{R})$ under the inner product $\langle f, g \rangle = \int_0^1 f(x)g(x) \, dx$.

  - ○ First, by applying Gram-Schmidt to the basis $\{1, x\}$, we can generate an orthogonal basis of $P_1(\mathbb{R})$ under this inner product: the result (after rescaling to clear denominators) is $\{1, 2x - 1\}$.

  - ○ Now, with $p_1 = 1$ and $p_2 = 2x - 1$, we can compute $\langle e^x, p_1 \rangle = e - 1$, $\langle e^x, p_2 \rangle = 3 - e$, $\langle p_1, p_1 \rangle = 1$, and $\langle p_2, p_2 \rangle = 1/3$.

  - ○ Then $\text{proj}_{P_2(\mathbb{R})}(e^x) = \dfrac{\langle e^x, p_1 \rangle}{\langle p_1, p_1 \rangle} p_1 + \dfrac{\langle e^x, p_2 \rangle}{\langle p_2, p_2 \rangle} p_2 = \boxed{(10 - 4e) + (18 - 6e)x \approx 0.873 + 1.690x}$.

  - ○ We can see from a plot (see above) of the orthogonal projection polynomial and $e^x$ on $[0, 1]$ that this line is indeed a very accurate approximation of $e^x$ on the interval $[0, 1]$.

## 3.3   Linear Transformations and Inner Products

- We will now study linear transformations $T : V \to W$ where $V$ and $W$ are (finite-dimensional) vector spaces, at least one of which is equipped with an inner product.

- Most of our results will use the fact that inner products are linear in their first coordinate: $\langle \mathbf{v}_1 + \mathbf{v}_2, \mathbf{w} \rangle = \langle \mathbf{v}_1, \mathbf{w} \rangle + \langle \mathbf{v}_2, \mathbf{w} \rangle$ and $\langle c\mathbf{v}, \mathbf{w} \rangle = c \langle \mathbf{v}, \mathbf{w} \rangle$.

### 3.3.1   Characterizations of Inner Products

- First, if we have an orthonormal basis for $W$, we can use it to compute the entries in the matrix associated to $T$.

- <u>Proposition</u> (Entries of Associated Matrix): Suppose $T : V \to W$ is a linear transformation of finite-dimensional vector spaces, and $W$ is an inner product space. If $\beta = \{\mathbf{v}_1, \ldots, \mathbf{v}_n\}$ is an ordered basis for $V$ and $\gamma = \{\mathbf{e}_1, \ldots, \mathbf{e}_m\}$ is an orthonormal ordered basis for $W$, then the $(i, j)$-entry of $[T]_\beta^\gamma$ is $\langle T(\mathbf{v}_j), \mathbf{e}_i \rangle$ for each $1 \leq i \leq m$ and $1 \leq j \leq n$.

  - ○ <u>Proof</u>: By our results on vector spaces with an orthonormal basis, we know that $T(\mathbf{v}_j) = \displaystyle\sum_{i=1}^{m} a_i \mathbf{e}_i$ where $a_i = \langle T(\mathbf{v}_j), \mathbf{e}_i \rangle$.

  - ○ But this is precisely the statement that the $i$th entry of the coefficient vector $[T(\mathbf{v}_j)]_\gamma$, which is in turn the $(i, j)$-entry of $[T]_\beta^\gamma$, is $\langle T(\mathbf{v}_j), \mathbf{e}_i \rangle$.

- <u>Example</u>: Let $T : P_2(\mathbb{R}) \to \mathbb{R}^3$ be defined by $T(p) = \langle p(0), p(1), p(2) \rangle$. Find the matrix associated to $T$ with respect to the bases $\beta = \{1 - 2x + x^2, 2x + x^2, 3 + x - x^2\}$ and $\gamma = \{\frac{1}{3} \langle 1, 2, 2 \rangle, \frac{1}{3} \langle 2, 1, -2 \rangle, \frac{1}{3} \langle -2, 2, -1 \rangle\}$ of $P_2(\mathbb{R})$ and $\mathbb{R}^3$ respectively.

- ○ Notice that $\gamma$ is an orthonormal basis of $\mathbb{R}^3$, so we need only compute the appropriate inner products.

- ○ Since $T(1 - 2x + x^2) = \langle 1, 0, 1 \rangle$, we compute $\langle 1, 0, 1 \rangle \cdot \frac{1}{3} \langle 1, 2, 2 \rangle = 1$, $\langle 1, 0, 1 \rangle \cdot \frac{1}{3} \langle 2, 1, -2 \rangle = 0$, and $\langle 1, 0, 1 \rangle \cdot \frac{1}{3} \langle -2, 2, -1 \rangle = -1$.

- ○ Next, since $T(2x + x^2) = \langle 0, 3, 8 \rangle$, we compute $\langle 0, 3, 8 \rangle \cdot \frac{1}{3} \langle 1, 2, 2 \rangle = 22/3$, $\langle 0, 3, 8 \rangle \cdot \frac{1}{3} \langle 2, 1, -2 \rangle = -13/3$, and $\langle 0, 3, 8 \rangle \cdot \frac{1}{3} \langle -2, 2, -1 \rangle = -2/3$.

- ○ Finally, since $T(3 + x - x^2) = \langle 3, 3, 1 \rangle$, we compute $\langle 3, 3, 1 \rangle \cdot \frac{1}{3} \langle 1, 2, 2 \rangle = 11/3$, $\langle 3, 3, 1 \rangle \cdot \frac{1}{3} \langle 2, 1, -2 \rangle = 7/3$, and $\langle 3, 3, 1 \rangle \cdot \frac{1}{3} \langle -2, 2, -1 \rangle = -1/3$.

- ○ Thus, $[T]_\beta^\gamma = \boxed{\begin{bmatrix} 1 & 22/3 & 11/3 \\ 0 & -13/3 & 7/3 \\ -1 & -2/3 & -1/3 \end{bmatrix}}$.

- As we noted above, for a fixed $\mathbf{w}$ in $V$, the function $T : V \to F$ defined by $T(\mathbf{v}) = \langle \mathbf{v}, \mathbf{w} \rangle$ is a linear transformation. In general, a linear transformation $T : V \to F$ from a vector space to its scalar field is called a <u>linear functional</u>.

  - ○ Perhaps surprisingly, when $V$ is a finite-dimensional inner product space, it turns out that every linear transformation $T : V \to F$ is of the form above.

- <u>Theorem</u> (Riesz Representation Theorem): If $V$ is a finite-dimensional inner product space with scalar field $F$, and $T : V \to F$ is linear, then there exists a unique vector $\mathbf{w}$ in $V$ such that $T(\mathbf{v}) = \langle \mathbf{v}, \mathbf{w} \rangle$.

  - ○ <u>Proof</u>: Let $\{\mathbf{e}_1, \ldots, \mathbf{e}_n\}$ be an orthonormal basis for $V$, and let $\mathbf{w} = \overline{T(\mathbf{e}_1)}\mathbf{e}_1 + \overline{T(\mathbf{e}_2)}\mathbf{e}_2 + \cdots + \overline{T(\mathbf{e}_n)}\mathbf{e}_n$.
  - ○ We claim that the linear transformation $R(\mathbf{v}) = \langle \mathbf{v}, \mathbf{w} \rangle$ is equal to $T$.
  - ○ For each $i$, $R(\mathbf{e}_i) = \left\langle \mathbf{e}_i, \overline{T(\mathbf{e}_1)}\mathbf{e}_1 + \overline{T(\mathbf{e}_2)}\mathbf{e}_2 + \cdots + \overline{T(\mathbf{e}_n)}\mathbf{e}_n \right\rangle = \left\langle \mathbf{e}_i, \overline{T(\mathbf{e}_i)}\mathbf{e}_i \right\rangle = \overline{\overline{T(\mathbf{e}_i)}} \langle \mathbf{e}_i, \mathbf{e}_i \rangle = T(\mathbf{e}_i)$.
  - ○ But now, since linear transformations are characterized by their values on a basis, the fact that $R(\mathbf{e}_i) = T(\mathbf{e}_i)$ for each $i$ implies that $R(\mathbf{v}) = T(\mathbf{v})$ for each $\mathbf{v}$ in $V$, as claimed.
  - ○ For uniqueness, if $\langle \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{w}' \rangle$ for all $\mathbf{v}$, then we would have $\langle \mathbf{v}, \mathbf{w} - \mathbf{w}' \rangle = 0$ for all $\mathbf{v}$. Setting $\mathbf{v} = \mathbf{w} - \mathbf{w}'$ yields $\langle \mathbf{w} - \mathbf{w}', \mathbf{w} - \mathbf{w}' \rangle = 0$, meaning that $\mathbf{w} - \mathbf{w}' = \mathbf{0}$ so that $\mathbf{w} = \mathbf{w}'$.

- We can also define a matrix associated to an inner product in a similar way to how we define a matrix associated to a linear transformation.

- <u>Definition</u>: If $V$ is a finite-dimensional inner product space with ordered bases $\beta$ and $\gamma$, we define the matrix $M_\beta^\gamma$ associated to the inner product on $V$ to have $(i, j)$-entry equal to $\langle \beta_j, \gamma_i \rangle$.

  - ○ Note that if $\beta = \gamma$ is an orthonormal basis for $V$, then the associated matrix will be the identity matrix.

- The main reason we use this definition is the following result, which is the inner-product analogue of the structure theorem for finite-dimensional vector spaces we proved earlier:

- <u>Theorem</u> (Inner Products and Dot Products): If $V$ is a finite-dimensional inner product space over $F$ with ordered bases $\beta$ and $\gamma$, and associated matrix $M_\beta^\gamma$, then for any vectors $\mathbf{v}$ and $\mathbf{w}$ in $V$, we have $\langle \mathbf{v}, \mathbf{w} \rangle = \left\langle M_\beta^\gamma [\mathbf{v}]_\beta, [\mathbf{w}]_\gamma \right\rangle$, where the second inner product is the standard inner product on $F^n$.

  - ○ If $F = \mathbb{R}$, then the standard inner product is the standard dot product $\mathbf{x} \cdot \mathbf{y}$, where if $F = \mathbb{C}$ then the standard inner product is the modified dot product $\mathbf{x} \cdot \overline{\mathbf{y}}$.
  - ○ This theorem therefore says that the inner product structure of $V$ is essentially the same as the (modified) dot product on $F^n$, up to a factor of the matrix $M_\beta^\gamma$. In particular, if we choose $\beta = \gamma$ to be an orthonormal basis, then the matrix factor is simply the identity matrix.

○ <u>Proof</u>: Suppose $\beta = \{\mathbf{v}_1, \ldots, \mathbf{v}_n\}$ and $\gamma = \{\mathbf{w}_1, \ldots, \mathbf{w}_n\}$, and take $\mathbf{v} = \sum_{i=1}^{n} a_i \mathbf{v}_i$ and $\mathbf{w} = \sum_{j=1}^{n} b_j \mathbf{w}_j$ for scalars $a_i$ and $b_j$.

○ Then by properties of inner products, $\langle \mathbf{v}, \mathbf{w} \rangle = \left\langle \sum_{i=1}^{n} a_i \mathbf{v}_i, \sum_{j=1}^{n} b_j \mathbf{w}_j \right\rangle = \sum_{i=1}^{n} \sum_{j=1}^{n} a_i \overline{b_j} \langle \mathbf{v}_i, \mathbf{w}_j \rangle$.

○ Also, since $[\mathbf{v}]_\beta = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$, we compute $M_\beta^\gamma [\mathbf{v}]_\beta = \begin{bmatrix} \sum_{i=1}^{n} a_i \langle \mathbf{v}_i, \mathbf{w}_1 \rangle \\ \vdots \\ \sum_{i=1}^{n} a_i \langle \mathbf{v}_i, \mathbf{w}_n \rangle \end{bmatrix}$.

○ So, since $[\mathbf{w}]_\gamma = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$, we see $\left\langle M_\beta^\gamma [\mathbf{v}]_\beta, [\mathbf{w}]_\gamma \right\rangle = \sum_{j=1}^{n} \overline{b_j} \left[ \sum_{i=1}^{n} a_i \langle \mathbf{v}_i, \mathbf{w}_j \rangle \right] = \sum_{j=1}^{n} \sum_{i=1}^{n} a_i \overline{b_j} \langle \mathbf{v}_i, \mathbf{w}_j \rangle$.

○ This expression is equal to the one we gave above, so we are done.

### 3.3.2 The Adjoint of a Linear Transformation

- As we have already noted, an inner product is linear in its first coordinate. Thus, if $T : V \to W$ is a linear transformation, the composition $\langle T(\mathbf{v}), \mathbf{w} \rangle$ will also be linear in its first coordinate, and as we have already seen, such compositions play an important role in the structure of an inner product.

- <u>Theorem</u> (Adjoint Transformations): Suppose that $V$ and $W$ are finite-dimensional inner product spaces with inner products $\langle \cdot, \cdot \rangle_V$ and $\langle \cdot, \cdot \rangle_W$ respectively, and that $T : V \to W$ is linear. Then there exists a unique function $T^* : W \to V$ such that $\langle T(\mathbf{v}), \mathbf{w} \rangle_W = \langle \mathbf{v}, T^*(\mathbf{w}) \rangle_V$ for all $\mathbf{v}$ in $V$ and $\mathbf{w}$ in $W$, and in fact $T^*$ is also a linear transformation.

    ○ <u>Proof</u>: For a fixed vector $\mathbf{w}$, the map $R_\mathbf{w}(\mathbf{v}) = \langle T(\mathbf{v}), \mathbf{w} \rangle_W$ is a linear transformation $R : V \to F$, since it is the composition of the linear transformations $\langle \cdot, \mathbf{w} \rangle_W$ and $T$.

    ○ Now, since $V$ is finite-dimensional, by the Riesz representation theorem, there exists a vector $\mathbf{w}'$ in $V$ such that $R_\mathbf{w}(\mathbf{v}) = \langle \mathbf{v}, \mathbf{w}' \rangle_V$ for every $\mathbf{v}$ in $V$.

    ○ Therefore, if we define $T^*(\mathbf{w})$ to be this vector $\mathbf{w}'$, we conclude that $\langle T(\mathbf{v}), \mathbf{w} \rangle_W = \langle \mathbf{v}, T^*(\mathbf{w}) \rangle_V$ for all $\mathbf{v}$ and $\mathbf{w}$ in $V$.

    ○ For uniqueness, if we had another function $T_0^*$ with $\langle T(\mathbf{v}), \mathbf{w} \rangle_W = \langle \mathbf{v}, T_0^*(\mathbf{w}) \rangle_V$, then we would have $\langle \mathbf{v}, T^*(\mathbf{w}) - T_0^*(\mathbf{w}) \rangle_V = \langle \mathbf{v}, T^*(\mathbf{w}) \rangle_V - \langle \mathbf{v}, T_0^*(\mathbf{w}) \rangle_V = \langle T(\mathbf{v}), \mathbf{w} \rangle_W - \langle T(\mathbf{v}), \mathbf{w} \rangle_W = 0$, so setting $\mathbf{v} = T^*(\mathbf{w}) - T_0^*(\mathbf{w})$ immediately gives $T^*(\mathbf{w}) = T_0^*(\mathbf{w})$.

    ○ Finally, to see that $T^*$ is linear, observe that

    $$\begin{aligned} \langle \mathbf{v}, T^*(\mathbf{w}_1 + \mathbf{w}_2) \rangle_V &= \langle T(\mathbf{v}), \mathbf{w}_1 + \mathbf{w}_2 \rangle_W = \langle T(\mathbf{v}), \mathbf{w}_1 \rangle_W + \langle T(\mathbf{v}), \mathbf{w}_2 \rangle_W = \langle \mathbf{v}, T^*(\mathbf{w}_1) + T^*(\mathbf{w}_2) \rangle_V \\ \langle \mathbf{v}, T^*(c\mathbf{w}) \rangle_V &= \langle T(\mathbf{v}), c\mathbf{w} \rangle_W = \overline{c} \langle T(\mathbf{v}), \mathbf{w} \rangle_W = \overline{c} \langle \mathbf{v}, T^*(\mathbf{w}) \rangle_V = \langle \mathbf{v}, cT^*(\mathbf{w}) \rangle_V \end{aligned}$$

    and so by the uniqueness of $T^*$ we must have $T^*(\mathbf{w}_1 + \mathbf{w}_2) = T^*(\mathbf{w}_1) + T^*(\mathbf{w}_2)$ and $T^*(c\mathbf{w}) = cT^*(\mathbf{w})$.

- <u>Definition</u>: If $V$ and $W$ are inner product spaces and $T : V \to W$ is linear, then, if it exists, the function $T^* : W \to V$ with the property that $\langle T(\mathbf{v}), \mathbf{w} \rangle_W = \langle \mathbf{v}, T^*(\mathbf{w}) \rangle_V$ for all $\mathbf{v}$ in $V$ and $\mathbf{w}$ in $W$ is called the <u>adjoint</u> of $T$.

    ○ The theorem above guarantees the existence and uniqueness of this map $T^*$ when the inner product spaces are finite-dimensional.

    ○ When $V$ and $W$ are infinite-dimensional, there may not exist such a map $T^*$.

    ○ It is not so easy to give an explicit counterexample, but if $V$ is the set of infinite sequences with only finitely many nonzero terms, and $T$ is the transformation that maps the standard basis vector $e_n$ to $e_1 + e_2 + \cdots + e_n$, then $T$ has no adjoint.

    ○ If $T$ does have an adjoint, however, the proof above shows that it is unique, and in any case, even when $T^*$ exists, it is not so obvious how to compute $T^*$ explicitly.

- <u>Example</u>: With $V = \mathbb{C}^2$ with the standard inner product and $T : V \to V$ is the linear transformation with $T(x, y) = (ax + by, cx + dy)$, find the adjoint map $T^*$.

○ Let $\{\mathbf{e}_1, \mathbf{e}_2\} = \{\langle 1, 0 \rangle, \langle 0, 1 \rangle\}$ be the standard orthonormal basis for $V$. Since a linear transformation is characterized by its values on a basis, it is enough to find $T^*(\mathbf{e}_1)$ and $T^*(\mathbf{e}_2)$.

○ By definition, $\langle (x, y), T^*(\mathbf{e}_1) \rangle = \langle T(x, y), \mathbf{e}_1 \rangle = \langle (ax + by, cx + dy), \mathbf{e}_1 \rangle = ax + by$.

○ Therefore, $T^*(\mathbf{e}_1)$ must be the vector $(\overline{a}, \overline{b})$, since that is the only vector with the property that $\langle (x, y), T^*(\mathbf{e}_1) \rangle = ax + by$ for arbitrary $x$ and $y$.

○ Similarly, $\langle (x, y), T^*(\mathbf{e}_2) \rangle = \langle T(x, y), \mathbf{e}_2 \rangle = \langle (ax + by, cx + dy), \mathbf{e}_2 \rangle = cx + dy$, so by the same argument as above, $T^*(\mathbf{e}_2)$ must be the vector $(\overline{c}, \overline{d})$.

○ Then we see $T^*(x, y) = xT^*(\mathbf{e}_1) + yT^*(\mathbf{e}_2) = \boxed{(\overline{a}x + \overline{c}y, \ \overline{b}x + \overline{d}y)}$.

○ Notice here that the matrix associated to $T$ is $[T]_\beta^\beta = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ while the matrix associated to $T^*$ is $[T^*]_\beta^\beta = \begin{bmatrix} \overline{a} & \overline{c} \\ \overline{b} & \overline{d} \end{bmatrix}$, the conjugate-transpose of the matrix associated to $T$.

- It seems, based on the example above, that $T^*$ is related to the conjugate-transpose of a matrix. In fact, this is true in general (which justifies our use of the same notation for both):

- <u>Proposition</u> (Adjoint Matrix): If $V$ and $W$ are finite-dimensional inner product spaces with orthonormal bases $\beta$ and $\gamma$ respectively, and $T : V \to W$ is linear, then the matrix $[T^*]_\gamma^\beta$ associated to the adjoint $T^*$ is $([T]_\beta^\gamma)^*$, the conjugate-transpose of the matrix associated to $T$.

  ○ <u>Proof</u>: Let $\beta = \{\mathbf{v}_1, \ldots, \mathbf{v}_n\}$ and $\gamma = \{\mathbf{w}_1, \ldots, \mathbf{w}_m\}$.

  ○ From our results on associated matrices in inner product spaces, the $(i, j)$-entry of $[T^*]_\gamma^\beta$ is $\langle T^*(\mathbf{w}_j), \mathbf{v}_i \rangle = \overline{\langle \mathbf{v}_i, T^*(\mathbf{w}_j) \rangle} = \overline{\langle T(\mathbf{v}_i), \mathbf{w}_j \rangle}$.

  ○ Now notice that the quantity $\overline{\langle T(\mathbf{v}_i), \mathbf{w}_j \rangle}$ is the complex conjugate of the $(j, i)$-entry of $[T]_\beta^\gamma$: thus, $[T^*]_\gamma^\beta = ([T]_\beta^\gamma)^*$ as claimed.

- Here are a few basic properties of the adjoint, which each follow via the uniqueness property:

  1. If $S : V \to W$ and $T : V \to W$ are linear, then $(S + T)^* = S^* + T^*$.

     ○ <u>Proof</u>: Observe that $\langle \mathbf{v}, (S + T)^*\mathbf{w} \rangle = \langle (S + T)\mathbf{v}, \mathbf{w} \rangle = \langle S(\mathbf{v}) + T(\mathbf{v}), \mathbf{w} \rangle = \langle \mathbf{v}, S^*(\mathbf{w}) + T^*(\mathbf{w}) \rangle$. Thus, since $(S + T)^*$ is unique, it must be equal to $S^* + T^*$.

  2. If $T : V \to W$ is linear, then $(cT)^* = \overline{c}\,T^*$.

     ○ <u>Proof</u>: Observe that $\langle \mathbf{v}, (cT)^*\mathbf{w} \rangle = \langle (cT)\mathbf{v}, \mathbf{w} \rangle = c \langle T(\mathbf{v}), \mathbf{w} \rangle = \langle \mathbf{v}, \overline{c}T^*(\mathbf{w}) \rangle$, so by uniqueness, $(cT)^* = \overline{c}T^*$.

  3. If $S : V \to W$ and $T : U \to V$ are linear, then $(ST)^* = T^*S^*$.

     ○ <u>Proof</u>: Observe that $\langle \mathbf{u}, (ST)^*\mathbf{w} \rangle = \langle (ST)\mathbf{u}, \mathbf{w} \rangle = \langle T(\mathbf{u}), S^*\mathbf{w} \rangle = \langle \mathbf{u}, T^*S^*(\mathbf{w}) \rangle$, so by uniqueness, $(ST)^* = T^*S^*$.

  4. If $T : V \to W$ is linear, then $(T^*)^* = T$.

     ○ <u>Proof</u>: Observe that $\langle T^*(\mathbf{v}), \mathbf{w} \rangle = \overline{\langle \mathbf{w}, T^*(\mathbf{v}) \rangle} = \overline{\langle T(\mathbf{w}), \mathbf{v} \rangle} = \langle \mathbf{v}, T(\mathbf{w}) \rangle$, so by uniqueness, $(T^*)^* = T$.

- As immediate corollaries, we can also deduce the analogous properties of conjugate-transpose matrices. (Of course, with the exception of the statement that $(AB)^* = B^*A^*$, these are all immediate from the definition.)

- We will mention one other useful property of adjoints:

- <u>Theorem</u> (Orthogonal Complements and Adjoints): Suppose $T : V \to W$ has an adjoint $T^* : W \to V$. Then $\ker(T)$ and $\text{im}(T^*)$ are orthogonal subspaces of $V$, and if $V$ is finite-dimensional they are orthogonal complements.

  ○ This result is a generalization of our earlier observation that the rowspace and nullspace of a real matrix $A$ are orthogonal.

  ○ <u>Proof</u>: Suppose $\mathbf{v} \in \ker(T)$ and $\mathbf{v}' \in \text{im}(T^*)$, so that $\mathbf{v}' = T^*\mathbf{w}$ for some $\mathbf{w} \in W$.

○ Then $\langle \mathbf{v}, \mathbf{v}' \rangle = \langle \mathbf{v}, T^* \mathbf{w} \rangle = \langle T\mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{0}, \mathbf{w} \rangle = 0$, and so $\mathbf{v}$ is orthogonal to $\mathbf{v}'$. This means $\ker(T)$ is orthogonal to $\mathrm{im}(T^*)$, as claimed. In particular, we have $\mathrm{im}(T^*) \subseteq [\ker(T)]^\perp$.

○ For the second statement, choose orthonormal bases $\beta$ of $V$ and $\gamma$ of $W$. Then for $A = [T]_\beta^\gamma$ we have $\dim(\mathrm{im}\, T) = \mathrm{rank}(A) = \mathrm{rank}(A^*) = \dim(\mathrm{im}\, T^*)$ since the rank of $A$ equals the rank of $A^*$ by our results on row and column spaces.

○ By our results on dimensions of orthogonal complements and the nullity-rank theorem, we have $\dim([\ker T]^\perp) = \dim(V) - \dim(\ker T) = \dim(\mathrm{im}\, T) = \dim(\mathrm{im}\, T^*)$.

○ But since $\mathrm{im}(T^*) \subseteq [\ker(T)]^\perp$ we must have equality since $V$ is finite-dimensional: thus $\mathrm{im}(T^*)$ is the orthogonal complement of $\ker(T)$, as claimed.

## 3.4 Applications of Inner Products

- In this section we discuss several practical applications of inner products.

### 3.4.1 Least-Squares Estimates

- A fundamental problem in applied mathematics and statistics is data fitting: finding a model that well approximates some set of experimental data. Problems of this type are ubiquitous in the physical sciences, social sciences, life sciences, and engineering.

  ○ A common example is that of finding a linear regression: a line $y = mx + b$ that best fits a set of 2-dimensional data points $\{(x_1, y_1), \ldots, (x_n, y_n)\}$ when plotted in the plane.

  ○ Of course, in many cases a linear model is not appropriate, and other types of models (polynomials, powers, exponential functions, logarithms, etc.) are needed instead.

  ○ The most common approach to such regression analysis is the method of "least squares", which minimizes the sum of the squared errors (the error being the difference between the model and the actual data).

- As we will discuss, many of these questions ultimately reduce to the following: if $A$ is an $m \times n$ matrix such that the matrix equation $A\mathbf{x} = \mathbf{c}$ has no solution, what vector $\hat{\mathbf{x}}$ is the closest approximation to a solution?

  ○ In other words, we are asking for the vector $\hat{\mathbf{x}}$ that minimizes the vector norm $||A\hat{\mathbf{x}} - \mathbf{c}||$.

  ○ Since the vectors of the form $A\hat{\mathbf{x}}$ are precisely those in the column space (i.e., image) of $A$, from our analysis of best approximations earlier we see that the vector $\mathbf{w} = A\hat{\mathbf{x}}$ will be the orthogonal projection of $\mathbf{c}$ into the column space of $A$.

  ○ Then, by orthogonal decomposition, we know that $\mathbf{w}^\perp = \mathbf{c} - A\hat{\mathbf{x}}$ is in the orthogonal complement of the column space of $A$.

  ○ By our theorem on orthogonal complements and adjoints, however, we know that the orthogonal complement of the column space of $A$ is the kernel of $A^*$.

  ○ Therefore, $\mathbf{w}^\perp$ is in $\ker(A^*)$, so $A^* \mathbf{w}^\perp = \mathbf{0}$.

  ○ Explicitly, this means $A^*(\mathbf{c} - A\hat{\mathbf{x}}) = \mathbf{0}$, which is to say, $A^* A\hat{\mathbf{x}} = A^* \mathbf{c}$: this is an explicit matrix system that we can solve for $\hat{\mathbf{x}}$.

- Definition: If $A$ is an $m \times n$ matrix with $m > n$, a least-squares solution to the matrix equation $A\mathbf{x} = \mathbf{c}$ is a vector $\hat{\mathbf{x}}$ satisfying $A^* A\hat{\mathbf{x}} = A^* \mathbf{c}$ (this equation is called the normal equation).

  ○ The system $A^* A\hat{\mathbf{x}} = A^* \mathbf{c}$ for $\hat{\mathbf{x}}$ is always consistent for any matrix $A$, although it is possible for there to be infinitely many solutions (a trivial case would be when $A$ is the zero matrix). Even in this case, the orthogonal projection $\mathbf{w} = A\hat{\mathbf{x}}$ onto the column space of $A$ will always be unique.

- In typical cases, the rank of $A$ is often equal to $n$. In this case, the matrix $A^* A$ will always be invertible, and there is a unique least-squares solution:

- Proposition (Least-Squares Solution): If $A$ is an $m \times n$ matrix and $\mathrm{rank}(A) = n$, then $A^* A$ is invertible and the unique least-squares solution to $A\mathbf{x} = \mathbf{c}$ is $\hat{\mathbf{x}} = (A^* A)^{-1} A^* \mathbf{c}$.

- We usually use this result in the situation where $A$ is real, in which case $A^* = A^T$.

- Proof: First, we show that the nullspace of $A^*A$ is the same as the nullspace of $A$. Clearly, if $A\mathbf{x} = \mathbf{0}$ then $(A^*A)\mathbf{x} = \mathbf{0}$. Conversely, if $A^*A\mathbf{x} = \mathbf{0}$ then $\langle A\mathbf{x}, A\mathbf{x} \rangle = \langle \mathbf{x}, A^*A\mathbf{x} \rangle = 0$, but by [I3] this means $A\mathbf{x} = \mathbf{0}$, and so $\mathbf{x}$ is in the nullspace of $A$.

- Now suppose $\text{rank}(A) = n$. Then since the dimension of the nullspace is the number of columns minus the rank, and $A^*A$ and $A$ both have $n$ columns, $\text{rank}(A^*A) = \text{rank}(A) = n$.

- But since $A^*A$ is an $n \times n$ matrix, this means $A^*A$ is invertible. The second statement then follows immediately upon left-multiplying $A\mathbf{x} = \mathbf{c}$ by $(A^*A)^{-1}$.

- Example: Find the least-squares solution to the inconsistent system $x + 2y = 3$, $2x + y = 4$, $x + y = 2$.

  - In this case, we have $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 1 \end{bmatrix}$ and $\mathbf{c} = \begin{bmatrix} 3 \\ 4 \\ 2 \end{bmatrix}$. Since $A$ clearly has rank 2, $A^*A$ will be invertible and there will be a unique least-squares solution.

  - We compute $A^*A = \begin{bmatrix} 6 & 5 \\ 5 & 6 \end{bmatrix}$, which is indeed invertible and has inverse $(A^*A)^{-1} = \dfrac{1}{11}\begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix}$.

  - The least-squares solution is therefore $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{c} = \boxed{\begin{bmatrix} 3 \\ -1 \end{bmatrix}}$.

  - In this case, we see $A\hat{\mathbf{x}} = \begin{bmatrix} 2 \\ 5 \\ 2 \end{bmatrix}$, so the error vector is $\mathbf{c} - A\hat{\mathbf{x}} = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$. Our analysis above indicates that this error vector has the smallest possible norm.

- We can apply these ideas to the problem of finding an optimal model for a set of data points.

  - For example, suppose that we wanted to find a linear model $y = mx + b$ that fits a set of data points $\{(x_1, y_1), \ldots, (x_n, y_n)\}$, in such a way as to minimize the sum of the squared errors $(y_1 - mx_1 - b)^2 + \cdots + (y_n - mx_n - b)^2$.

  - If the data points happened to fit exactly on a line, then we would be seeking the solution to the system $y_1 = mx_1 + b$, $y_2 = mx_2 + b$, ... , $y_n = mx_n + b$.

  - In matrix form, this is the system $A\mathbf{x} = \mathbf{c}$ where $A = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} b \\ m \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$.
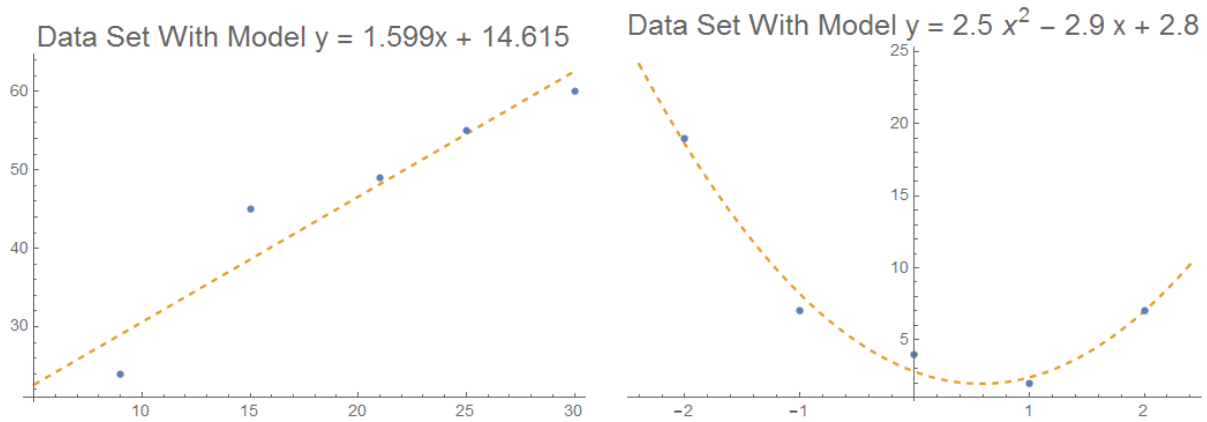
  - Of course, due to experimental errors and other random noise, it is unlikely for the data points to fit the model exactly. Instead, the least-squares estimate $\hat{\mathbf{x}}$ will provide the values of $m$ and $b$ that minimize the sum of the squared errors.

  - In a similar way, to find a quadratic model $y = ax^2 + bx + c$ for a data set $\{(x_1, y_1), \ldots, (x_n, y_n)\}$, we would use the least-squares estimate for $A\mathbf{x} = \mathbf{c}$, with $A = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} c \\ b \\ a \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$.

  - In general, to find a least-squares model of the form $y = a_1 f_1(x) + \cdots + a_m f_m(x)$ for a data set $\{(x_1, y_1), \ldots, (x_n, y_n)\}$, we would want the least-squares estimate for the system $A\mathbf{x} = \mathbf{c}$, with $A = \begin{bmatrix} f_1(x_1) & \cdots & f_m(x_1) \\ \vdots & \ddots & \vdots \\ f_1(x_n) & \cdots & f_m(x_n) \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$.

- Example: Use least-squares estimation to find the line $y = mx + b$ that is the best model for the data points $\{(9, 24), (15, 45), (21, 49), (25, 55), (30, 60)\}$.

◦ We seek the least-squares solution for $A\mathbf{x} = \mathbf{c}$, where $A = \begin{bmatrix} 1 & 9 \\ 1 & 15 \\ 1 & 21 \\ 1 & 25 \\ 1 & 30 \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} b \\ m \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} 24 \\ 45 \\ 49 \\ 55 \\ 60 \end{bmatrix}$.

◦ We compute $A^*A = \begin{bmatrix} 5 & 100 \\ 100 & 2272 \end{bmatrix}$, so the least-squares solution is $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{c} \approx \begin{bmatrix} 14.615 \\ 1.599 \end{bmatrix}$.

◦ Thus, to three decimal places, the desired line is $y = \boxed{1.599x + 14.615}$. From a plot, we can see that this line is fairly close to all of the data points:



Data Set With Model y = 1.599x + 14.615



Data Set With Model y = 2.5 $x^2$ − 2.9 x + 2.8

• Example: Use least-squares estimation to find the quadratic function $y = ax^2 + bx + c$ best modeling the data points $\{(-2, 19), (-1, 7), (0, 4), (1, 2), (2, 7)\}$.

◦ We seek the least-squares solution for $A\mathbf{x} = \mathbf{c}$, with $A = \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} c \\ b \\ a \end{bmatrix}$, $\mathbf{c} = \begin{bmatrix} 19 \\ 7 \\ 4 \\ 2 \\ 7 \end{bmatrix}$.

◦ We compute $A^*A = \begin{bmatrix} 5 & 0 & 10 \\ 0 & 10 & 0 \\ 10 & 0 & 34 \end{bmatrix}$, so the least-squares solution is $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{c} = \begin{bmatrix} 2.8 \\ -2.9 \\ 2.5 \end{bmatrix}$.
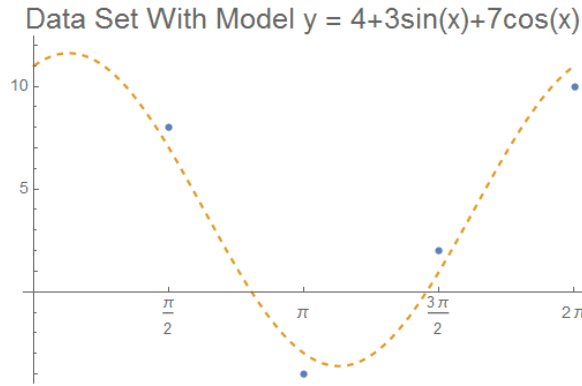
◦ Thus, the desired quadratic polynomial is $y = \boxed{-2.5x^2 - 2.9x + 2.8}$. From a plot (see above), we can see that this quadratic function is fairly close to all of the data points.

• Example: Use least-squares estimation to find the trigonometric function $y = a + b\sin(x) + c\cos(x)$ best modeling the data points $\{(\pi/2, 8), (\pi, -4), (3\pi/2, 2), (2\pi, 10)\}$.

◦ We seek the least-squares solution for $A\mathbf{x} = \mathbf{c}$, with $A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & -1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$, $\mathbf{c} = \begin{bmatrix} 8 \\ -4 \\ 2 \\ 10 \end{bmatrix}$.

◦ We compute $A^*A = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$, so the least-squares solution is $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{c} = \begin{bmatrix} 4 \\ 3 \\ 7 \end{bmatrix}$.

◦ Thus, the desired function is $y = \boxed{4 + 3\sin(x) + 7\cos(x)}$. In this case, the model predicts the points $\{(\pi/2, 7), (\pi, -3), (3\pi/2, 1), (2\pi, 11)\}$, so it is a good fit to the original data:

Data Set With Model y = 4+3sin(x)+7cos(x)

- Our results on least squares also yield a method for writing down the matrix associated to orthogonal projection onto a subspace $W$ of $F^n$ with respect to the standard basis:

- <u>Corollary</u> (Associated Matrices for Projections): Suppose $F = \mathbb{R}$ or $\mathbb{C}$ and $W$ is a subspace of $F^n$ with basis $\{\mathbf{v}_1, \ldots, \mathbf{v}_k\}$. If $\beta$ is the standard basis for $F^n$ and $T : F^n \to F^n$ represents orthogonal projection onto the subspace $S$, then the associated matrix $[T]_\beta^\beta = A(A^*A)^{-1}A^*$, where $A$ is the $n \times k$ matrix whose columns are the vectors $\mathbf{v}_i$.

  ○ <u>Proof</u>: By our earlier results, for any column vector $\mathbf{c}$, the unique least-squares solution to $A\mathbf{x} = \mathbf{c}$ is $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{c}$, and the vector $A\hat{\mathbf{x}} = A(A^TA)^{-1}A^T\mathbf{c}$ represents the projection of $\mathbf{c}$ into the column space of $A$.

  ○ But the column space of the matrix $A$ is precisely $W$, by definition, so the associated matrix $[T]_\beta^\beta$ is precisely $A(A^*A)^{-1}A^*$, as claimed.

- <u>Example</u>: Find the matrix $M$ (with respect to the standard basis of $\mathbb{R}^4$) associated to orthogonal projection onto the subspace $W$ spanned by $\{(1,1,0,0),(-1,1,1,1)\}$ inside $\mathbb{R}^4$.

  ○ By the corollary above, the associated matrix is $A(A^*A)^{-1}A^*$ where $A = \begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}$.

  ○ We compute $(A^*A)^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/4 \end{bmatrix}$ and so $M = A(A^*A)^{-1}A^* = \dfrac{1}{4}\begin{bmatrix} 3 & 1 & -1 & -1 \\ 1 & 3 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{bmatrix}$.

  ○ To verify that this really is the correct matrix, we can check that the column space of $M$ is equal to $W$ (which it is), and that $M$ sends each vector in $W$ to itself (which it does).

- In addition to showing up in least-squares methods (which are extremely important in statistics and the experimental sciences), orthogonal projections also have many applications in computer graphics, coding theory, and machine learning.

  ○ To illustrate the idea very generally, if we have a set of data that is high-dimensional (i.e., lies inside $\mathbb{R}^n$ where $n$ is very large) that has a lot of underlying structure, it is often the case that projecting onto a much smaller-dimensional subspace will not lose very much information. Storing the projection of the data then requires much less information, which is the central idea of data compression.

  ○ To minimize the loss of information when compressing data, one may use tools such as principal-component analysis, which provide ways to calculate subspaces that carry as much of the information from the original data set as possible.

### 3.4.2 Fourier Series

- Another extremely useful application of the general theory of orthonormal bases is that of Fourier series.

  ○ Fourier analysis, broadly speaking, studies the problem of approximating a function on an interval by trigonometric functions. This problem is very similar to the question, studied in calculus, of approximating a function by a polynomial (the typical method is to use Taylor polynomials, although as we have already discussed, least-squares estimates provide another potential avenue).

  ○ Fourier series have a tremendously wide variety of applications, ranging from to solving partial differential equations (in particular, the famous wave equation and heat equation), studying acoustics and optics (decomposing an acoustic or optical waveform into simpler waves of particular frequencies), electical engineering, and quantum mechanics.

- Although a full discussion of Fourier series belongs more properly to analysis, we can still provide an overview.

  ○ A typical scenario in Fourier analysis is to approximate a continuous function on $[0, 2\pi]$ using a trigonometric polynomial: a function that is a polynomial in $\sin(x)$ and $\cos(x)$.

  ○ Using trigonometric identities, this question is equivalent to approximating a function $f(x)$ by a (finite) <u>Fourier series</u> of the form $s(x) = a_0 + b_1 \cos(x) + b_2 \cos(2x) + \cdots + b_k \cos(kx) + c_1 \sin(x) + c_2 \sin(2x) + \cdots + c_k \sin(kx)$.

  ○ Notice that, in the expression above, $s(0) = s(2\pi)$ since each function in the sum has period $2\pi$. Thus, we can only realistically hope to get close approximations to functions satisfying $f(0) = f(2\pi)$.

- Let $V$ be the vector space of continuous, real-valued functions on the interval $[0, 2\pi]$ having equal values at 0 and $2\pi$, and define an inner product on $V$ via $\langle f, g \rangle = \int_0^{2\pi} f(x)g(x)\, dx$.

- <u>Proposition</u>: The functions $\{\varphi_0, \varphi_1, \varphi_2, \dots\}$ are an orthonormal set on $V$, where $\varphi_0(x) = \dfrac{1}{\sqrt{2\pi}}$, and $\varphi_{2k-1}(x) = \dfrac{1}{\sqrt{\pi}} \cos(kx)$ and $\varphi_{2k}(x) = \dfrac{1}{\sqrt{\pi}} \sin(kx)$ for each $k \geq 1$.

  ○ <u>Proof</u>: Using the product-to-sum identities, such as $\sin(ax)\sin(bx) = \dfrac{1}{2}\left[\cos(a-b)x - \cos(a+b)x\right]$, it is a straightforward exercise in integration to verify that $\langle \varphi_i, \varphi_j \rangle = 0$ for each $i \neq j$.

  ○ Furthermore, we have $\langle \varphi_0, \varphi_0 \rangle = \dfrac{1}{2\pi} \int_0^{2\pi} 1\, dx = 1$, $\langle \varphi_{2k-1}, \varphi_{2k-1} \rangle = \dfrac{1}{\pi} \int_0^{2\pi} \cos^2(kx)\, dx = 1$, and $\langle \varphi_{2k}, \varphi_{2k} \rangle = \dfrac{1}{\pi} \int_0^{2\pi} \sin^2(kx)\, dx = 1$. Thus, the set is orthonormal.

- If it were the case that $S = \{\varphi_0, \varphi_1, \varphi_2, \dots\}$ were an orthonormal basis for $V$, then, given any other function $f(x)$ in $V$, we could write $f$ as a linear combination of functions in $\{\varphi_0, \varphi_1, \varphi_2, \dots\}$, where we can compute the appropriate coefficients using the inner product on $V$.

  ○ Unfortunately, $S$ does not span $V$: we cannot, for example, write the function $g(x) = \displaystyle\sum_{n=1}^{\infty} \dfrac{1}{2^n} \sin(nx)$ as a finite linear combination of $\{\varphi_0, \varphi_1, \varphi_2, \dots\}$, since doing so would require each of the infinitely many terms in the sum.

  ○ Ultimately, the problem, as exemplified by the function $g(x)$ above, is that the definition of "basis" only allows us to write down finite linear combinations.

  ○ On the other hand, the finite sums $\displaystyle\sum_{j=0}^{k} a_j \varphi_j(x)$ for $k \geq 0$, where $a_j = \langle f, \varphi_j \rangle$, will represent the best approximation to $f(x)$ inside the subspace of $V$ spanned by $\{\varphi_0, \varphi_1, \dots, \varphi_k\}$. Furthermore, as we increase $k$, we are taking approximations to $f$ that lie inside larger and larger subspaces of $V$, so as we take $k \to \infty$, these partial sums will yield better and better approximations to $f$.

  ○ Provided that $f$ is a sufficiently nice function, it can be proven that in the limit, our formulas for the coefficients do give a formula for $f(x)$ as an infinite sum:

- <u>Theorem</u> (Fourier Series): Let $f(x)$ be a twice-differentiable function on $[0, 2\pi]$ satisfying $f(0) = f(2\pi)$, and define the <u>Fourier coefficients</u> of $f$ as $a_j = \langle f, \varphi_j \rangle = \int_0^{2\pi} f(x)\varphi_j(x)\,dx$, for the trigonometric functions $\varphi_j(x)$ defined above. Then $f(x)$ is equal to its <u>Fourier series</u> $\sum_{j=0}^{\infty} a_j \varphi_j(x)$ for every $x$ in $[0, 2\pi]$.

  ○ This result can be interpreted as a "limiting version" of the theorem we stated earlier giving the coefficients for the linear combination of a vector in terms of an orthonormal basis: it gives an explicit way to write the function $f(x)$ as an "infinite linear combination" of the orthonormal basis elements $\{\varphi_0, \varphi_1, \varphi_2, \dots\}$.

- <u>Example</u>: Compute the Fourier coefficients and Fourier series for $f(x) = (x - \pi)^2$ on the interval $[0, 2\pi]$, and compare the partial sums of the Fourier series to the original function.
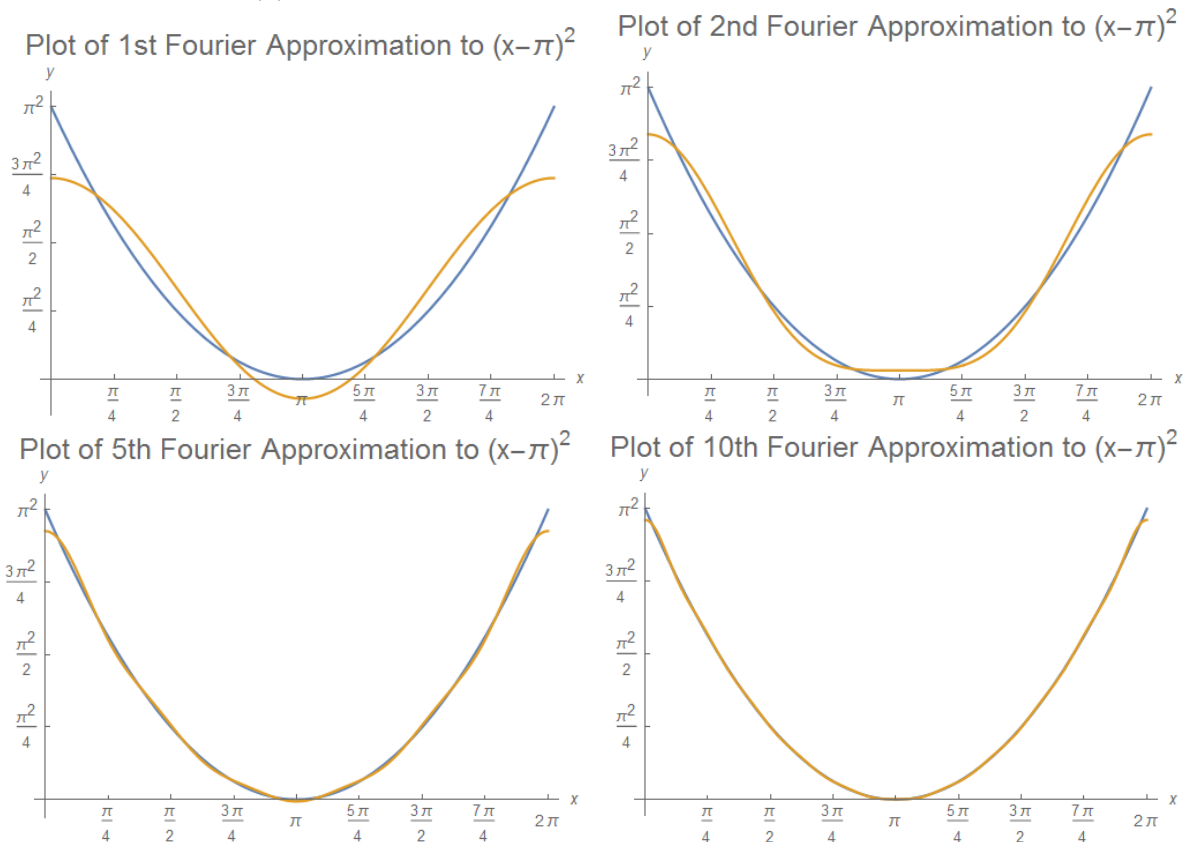
  ○ First, we have $a_0 = \int_0^{2\pi} f(x) \dfrac{1}{\sqrt{2\pi}}\,dx = \dfrac{1}{\sqrt{18}}\pi^{5/2}$.

  ○ For $k$ odd, after integrating by parts twice, we have $a_{2k-1} = \int_0^{2\pi} f(x) \dfrac{1}{\sqrt{\pi}}\cos(kx)\,dx = \dfrac{4\sqrt{\pi}}{k^2}$.

  ○ For $k$ even, in a similar manner we see $a_{2k} = \int_0^{2\pi} f(x) \dfrac{1}{\sqrt{\pi}}\sin(kx)\,dx = 0$.

  ○ Therefore, the Fourier series for $f(x)$ is $\boxed{\dfrac{1}{6}\pi^2 + \sum_{k=1}^{\infty} \dfrac{4}{k^2}\cos(kx)}$.

  ○ Here are some plots of the partial sums (up to the term involving $\cos(nx)$) of the Fourier series along with $f$. As is clearly visible from the graphs, the partial sums give increasingly close approximations to the original function $f(x)$ as we sum more terms:



Plot of 1st Fourier Approximation to $(x-\pi)^2$



Plot of 2nd Fourier Approximation to $(x-\pi)^2$



Plot of 5th Fourier Approximation to $(x-\pi)^2$



Plot of 10th Fourier Approximation to $(x-\pi)^2$

Well, you're at the end of my handout. Hope it was helpful.