**308 Review**

Much of what follows in this "review" will look very new. Some of it will look new because I am phrasing what you learned in 308 in a different way. Some of it will look new because it was not taught in your 308 class.

**Definition 0.1.** A **vector space** is a triple, $(V, +, \mathbb{F})$, where

1. $V$ is a set, a collection of objects. Objects in $V$, written
$$\vec{v} \in V,$$
   are called **vectors.**

2. $+$ is an operation which takes two objects in $V$ and produces another object in $V$. We write this
$$+ : V \times V \to V.$$
   We call this operation **addition.** This means $V$ must be **closed under addition.**

3. $\mathbb{F}$ is called a **scalar field**. In this class, $\mathbb{F}$ will either be $\mathbb{C}$, the complex numbers, or $\mathbb{R}$, the real numbers. Recall that we define
$$\mathbb{C} = \{a + ib \mid a \in \mathbb{R}, b \in \mathbb{R}\}$$
   Addition of complex numbers follows the rule,
$$(a + ib) + (c + id) = (a + c) + i(b + d).$$
   We can multiply complex numbers, as well, by the rule,
$$(a + ib) \cdot (c + id) = ac + ibc + iad + i^2 bd = (ac - bd) + i(bc + ad).$$

4. The set, $V$, must be **closed under scalar multiplication**. That is, for all vectors, $\vec{v} \in V$, and all $z \in \mathbb{F}$,
$$z\vec{v} \in V$$

5. There must be a special vector, written $\vec{0} \in V$, which has the property that for any $\vec{v} \in V$,
$$\vec{v} + \vec{0} = \vec{v}.$$

6. For each $\vec{v} \in V$, there must exist a vector, $-\vec{v}$ such that
$$\vec{v} + (-\vec{v}) = \vec{0}.$$

7. The operation, $+$ and scalar multiplication must satisfy the following properties for all $\vec{x}, \vec{y}, \vec{z} \in V$ and all $a, b \in \mathcal{F}$

- $(\vec{x} + \vec{y}) = (\vec{y} + \vec{x})$
- $(\vec{x} + \vec{y}) + \vec{z} = \vec{x} + (\vec{y} + \vec{z})$
- $a(\vec{x} + \vec{y}) = a\vec{x} + a\vec{y}$
- $a(b\vec{x}) = (ab)\vec{x}$
- $(a + b)\vec{x} = a\vec{x} + b\vec{x}$

**Remark 0.1.** *Note, that "closed under addition" and "closed under scalar multiplication"* **are properties of a set.** *The set is closed under the operation of addition or scalar multiplication.*

**Definition 0.2.** Given a vector space, $(V, +, \mathbb{F})$, a **subspace of** $(V, +, \mathbb{F})$ is a vector space, $(V', +, \mathbb{F})$ such that

$$V' \subset V$$

The operation, $+$, and the scalar field, $\mathbb{F}$ (think $\mathbb{R}$ or $\mathbb{C}$), must be the same.

In 308, the only vector spaces you saw were $(\mathbb{R}^n, +, \mathbb{R})$ and their subspaces. In this class we will be dealing with more general vector spaces.

**Definition 0.3.** Let $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}$, be a collection of vectors in a vector space, $(V, +, \mathbb{F})$. We say that a vector $\vec{v} \in V$ is a **linear combination** of the vectors $\vec{v}_1, \vec{v}_2, ..., \vec{v}_n$, if there exist scalars, $c_1, c_2, ..., c_n \in \mathbb{F}$ such that

$$\vec{v} = c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_n\vec{v}_n$$

If all $c_i = 0$, then $\vec{v} = \vec{0}$ and we call this the **trivial combination.** If there exists a $c_i \neq 0$, then we say that this is a **non-trivial combination.**

**Definition 0.4.** The **span of a collection of vectors**, $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}$, is defined as the following set,

$$span(\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}) = \{\vec{v}| \quad \vec{v} \text{ is a linear combination of } \{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}\}.$$

**The span of any collection of vectors is a subspace.**

**Definition 0.5.** A collection of vectors, $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}$, is called **linearly independent** if the only linear combination such that

$$c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_n\vec{v}_n = \vec{0}$$

is the trivial combination. The collection is called **linearly dependent** if there exists a non-trivial combination such that

$$c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_n\vec{v}_n = \vec{0}.$$

**Definition 0.6.** Given a vector space, $(V, +, \mathbb{F})$, a collection of vectors, $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\} \subset V$ , is called a **basis for** $(V, +, \mathbb{F})$ if it satisfies the following two properties,

1. $span(\vec{v}_1, \vec{v}_2, ..., \vec{v}_n) = V$.

2. The collection, $\{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}$, is linearly independent.

The vectors in a basis are called **basis vectors.**

Once we have a basis for $V$, it is true that every vector, $v \in V$, can be written as a unique linear combination of basis vectors.

**Definition 0.7.** Let $\mathcal{B} = \{\vec{v}_1, \vec{v}_2, ..., \vec{v}_n\}$ be a basis for a vector space $(V, +, \mathbb{F})$. For each $\vec{v} \in V$, we define the **coordinate vector with respect to the basis** $\mathcal{B}$ as

$$\vec{v}_{\mathcal{B}} = \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix}$$

where $c_1, c_2, ..., c_n$ are the coefficients of the unique linear combination such that

$$\vec{v} = c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_n\vec{v}_n.$$

In 308, when we were dealing with $\mathbb{R}^n$, we had always chosen the standard basis. However, sometimes there is a more convenient basis.

**Definition 0.8.** A **real-valued** $(n \times m)-$**matrix** is a rectangular array of $n$ rows and $m$ columns, where all the entries are real numbers. We denote the collection of all such matrices by,

$$M_{n \times m}(\mathbb{R}) = \{A| \quad A \text{ is a real-valued } (n \times m) - \text{ matrix}\}.$$

A **complex-valued** $(n \times m)-$**matrix** is a rectangular array of $n$ rows and $m$ columns, where all the entries are complex numbers. We denote the collection of all such matrices by,

$$M_{n \times m}(\mathbb{C}) = \{A| \quad A \text{ is a complex-valued } (n \times m) - \text{ matrix}\}.$$

Note that because $\mathbb{R} \subset \mathbb{C}$, we also have that $M_{n \times m}(\mathbb{R}) \subset M_{n \times m}(\mathbb{C})$.

We shall think of matrices $A \in M_{n \times m}(\mathbb{C})$ as rectangular arrays of $m$ columns, each of which is a coordinate vector in $\mathbb{C}^n$. When we wish to emphasize this perspective, we shall write

$$A = \begin{pmatrix} \vec{v}_1 & \vec{v}_1 & \cdots & \vec{v}_m \end{pmatrix}$$

where each $\vec{v}_i \in \mathbb{C}^n$

**Definition 0.9.** Let $A \in M_{n \times m}(\mathbb{C})$ be written as

$$A = \begin{pmatrix} \vec{v}_1 & \vec{v}_1 & \cdots & \vec{v}_m \end{pmatrix}$$

where each $\vec{v}_i \in \mathbb{C}^n$. Let $\vec{x} \in \mathbb{C}^m$ be written as the coordinate vector, $\vec{x} = \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{pmatrix}$. We define

**matrix-vector multiplication** according to the following rule,

$$A\vec{x} = c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_m\vec{v}_m$$

Note that $A\vec{x} \in \mathbb{C}^n$.

**Definition 0.10.** Let $A \in M_{n \times m}(\mathbb{C})$ and $B \in M_{m \times r}(\mathbb{C})$ where $B$ is written as

$$B = \begin{pmatrix} \vec{b}_1 & \vec{b}_1 & \cdots & \vec{b}_r \end{pmatrix}$$

for $\vec{b}_i \in \mathbb{C}^m$. We define **matrix-matrix multiplication** according to the following rule,

$$AB = \begin{pmatrix} A\vec{b}_1 & A\vec{b}_1 & \cdots & A\vec{b}_m \end{pmatrix}$$

Note that $AB \in M_{n \times r}(\mathbb{C})$.

Matrix multiplication is not always commutative. That is, even when it is defined, it may very well be that $AB \neq BA$.

**Definition 0.11.** Consider $\mathbb{R}^n$. Let

$$\mathcal{B}_1 = \{\vec{b}_i\}_{i=1}^n, \quad \mathcal{B}_2 = \{\vec{v}_i\}_{i=1}^n$$

be two different bases for $\mathbb{R}^n$. The **change of basis matrix from $\mathcal{B}_1$ to $\mathcal{B}_2$** is the unique $(n \times n)-$matrix, $S$, such that for all $\vec{x} \in \mathbb{R}^n$

$$\vec{x}_{\mathcal{B}_2} = S\vec{x}_{\mathcal{B}_1}.$$

Note that the change of basis matrix, $S$, is always invertible. If we write out $S$ as a matrix of column vectors,

$$S = \begin{pmatrix} \vec{s}_1 & \vec{s}_2 & \cdots & \vec{s}_n \end{pmatrix}$$

then $\vec{s}_i = (\vec{b}_i)_{\mathcal{B}_2}$. In other words, $\vec{s}_i = \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix}$ where the constants $c_i$ are the constants of the

unique linear combination such that

$$\vec{b}_i = c_1\vec{v}_1 + c_2\vec{v}_2 + ... + c_n\vec{v}_n.$$

**Reflection Questions**

1. Show that $(M_{n \times m}(\mathbb{C}), +, \mathbb{C})$ is a vector space where addition is defined by component-wise addition and scalar multiplication acts by scaling each component. You must use the definition.

2. Show $((M_{n \times m}(\mathbb{R}), +, \mathbb{R}))$ is a subspace of $(M_{n \times m}(\mathbb{C}), +, \mathbb{C})$. Show $((M_{n \times m}(\mathbb{R}), +, \mathbb{C}))$ is a not subspace of $(M_{n \times m}(\mathbb{C}), +, \mathbb{C})$ because it is not a vector space.

3. Let
$$C^1([a,b], \mathbb{R}) = \{f : [a,b] \to \mathbb{R} \mid f' \text{ is continuous}\}.$$
   Show that $(C^1([a,b], \mathbb{R}), +, \mathbb{R})$ is a vector space.

4. Is $\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$ in the span of the collection $\{ \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -3 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix} \}$?

5. Is the collection $\{ \begin{pmatrix} 1 \\ 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 10 \\ -3 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ -6 \\ 1 \end{pmatrix} \}$ linearly independent? A basis?

6. Multiply $\begin{pmatrix} 1 & 2 & 7 \\ 2 & 0 & 1 \\ 0 & 3 & -3 \end{pmatrix} \begin{pmatrix} 10 \\ 0 \\ -2 \end{pmatrix}$.

7. Let $\mathcal{B}_1 = \{ \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 3 \end{pmatrix} \}$ and $\mathcal{B}_2 = \{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 3 \\ 3 \end{pmatrix} \}$. Find the change of basis matrix which takes vectors written in $\mathcal{B}_1$ and expresses them in $\mathcal{B}_2$.

8. For $\mathcal{B}_1$ and $\mathcal{B}_2$ as above, find the change of basis matrix which takes vectors written in $\mathcal{B}_2$ and expresses them in $\mathcal{B}_1$.

# 1  Matrices and Eigenstuff

**Definition 1.1.** Let $(X, +_X, \mathbb{F})$ and $(Y, +_Y, \mathbb{F})$ be vector spaces. A map, $f : X \to Y$ is called **linear** if and only if for all $x_1, x_2 \in X$ and all $c \in \mathbb{F}$,

$$f(x_1 +_X cx_2) = f(x_1) +_Y cf(x_2).$$

### Examples

1. The function $f : \mathbb{R} \to \mathbb{R}$ given by the algebraic rule $f(x) = 2x$ is a linear function.

2. The function $f : \mathbb{R}^2 \to \mathbb{R}$ given by the algebraic rule $f(\vec{x}) = \vec{x} \cdot \vec{y_0}$ is a linear function.

3. Given any interval, $[a, b]$, the operation of taking the anti-derivative is a linear operation. That is, the map $int$ on functions $f$ defined by

$$int(f) = \int_a^b f(x)dx$$

is a linear map from integrable functions on $[a, b]$ to $\mathbb{R}$.

4. The operation of taking the derivative is also a linear operation. That is, the map, $\frac{d}{dt}$ defined by

$$\frac{d}{dt}(f) = f'$$

is a linear map from the vector space of differentiable scalar-valued functions to the vector space of scalar-valued functions.

In 308, you studied Linear Transformations, functions $T : \mathbb{R}^n \to \mathbb{R}^m$ which were linear. In particular you saw the following important idea:

**Theorem 1.1.** *(Representation of Linear Transformations) For every $(n \times m)-$ matrix, $A$, we can define a linear transformation, $T : \mathbb{R}^m \to \mathbb{R}^n$, by the rule*

$$T(\vec{x}) = A\vec{x}$$

*Conversely, for every linear transformation, $T : \mathbb{R}^m \to \mathbb{R}^n$, there is an $(n \times m)-$ matrix, $A$, such that*

$$T(\vec{x}) = A\vec{x}$$

*for all $\vec{x} \in \mathbb{R}^m$.*

*Because of this, we say that **the matrix, $A$, represents the linear transformation, $T$, and $T$ acts by multiplication by the matrix $A$.***

**Question 1.1:.** Given a linear transformation, $T : \mathbb{R}^m \to \mathbb{R}^n$, how do we find the matrix by which represents it? I.e., how do we find the matrix $A$ such that $T(\vec{x}) = A\vec{x}$?

**Answer:** Linear Transformations (and linear maps in general) are nice because the respect vector space operations like taking linear combinations. For example, if we let $\mathcal{S} = \{\vec{e}_i\}_{i=1}^m$ be the standard basis, and $\vec{x} = \sum_{i=1}^m x_i \vec{e}_i$, then

$$T(\vec{x}) = T(\sum_{i=1}^m x_i \vec{e}_i) = \sum_{i=1}^m x_i T(\vec{e}_i).$$

This means that linear transformations are completely determined by how they act upon bases. If we know how a linear transformation acts upon a basis, we know how it acts on every vector in the span of that basis. Thus, if we let $A$ be the matrix which represents $T$, we see that

$$
\begin{aligned}
A &= AI_m &= A[\vec{e}_1 \cdots \vec{e}_m] \\
&&= [A\vec{e}_1 \cdots A\vec{e}_m] \\
&&= [T(\vec{e}_1) \cdots T(\vec{e}_m)].
\end{aligned}
$$

**Note:** There are infinitely many bases. The matrix we found above represents $T$ in the standard basis. $T(\vec{x}) = A\vec{x}$ **only if $\vec{x}$ is written as a coordinate vector with respect to the standard basis.**

That is, for our vector $\vec{x} = \sum_{i=1}^m x_i \vec{e}_i$, we write $\vec{x}$ as a coordinate vector with respect to the standard basis as

$$\vec{x}_{\mathcal{S}} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$

Then, we have that,

$$T(\vec{x}) = T(\sum_{i=1}^m x_i \vec{e}_i) = \sum_{i=1}^m x_i T(\vec{e}_i),$$

and,

$$\sum_{i=1}^m x_i T(\vec{e}_i) = [T(\vec{e}_1) \cdots T(\vec{e}_m)] \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} = A\vec{x}_{\mathcal{S}}.$$

If we had written $\vec{x}$ as a coordinate vector with respect to a different basis, say, $\vec{x} = \sum_{i=1}^m c_i \vec{b}_i$, then

$$T(\vec{x}) = T(\sum_{i=1}^m c_i \vec{b}_i) = \sum_{i=1}^m c_i T(\vec{b}_i),$$

BUT,

$$\sum_{i=1}^{m} c_i T(\vec{b}_i) \neq [T(\vec{e}_1) \cdots T(\vec{e}_m)] \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix} = A\vec{x}_\mathcal{B}.$$

Rather, we would have

$$\sum_{i=1}^{m} c_i T(\vec{b}_i) = [T(\vec{b}_1) \cdots T(\vec{b}_m)] \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix} = B\vec{x}_\mathcal{B}.$$

**Question 1.2:.** How do we find a representation of a given a linear transformation, $T : \mathbb{R}^m \to \mathbb{R}^n$, in a given basis, $\mathcal{B} = \{\vec{b}_i\}_{i=1}^{m}$? I.e., how do we find the matrix $M$ such that $T(\vec{x}_\mathcal{B}) = M\vec{x}_\mathcal{B}$

**Answer:** This question is about change of bases. So, let's assume that the vectors of the basis $\mathcal{B} = \{\vec{b}_i\}_{i=1}^{m}$ are written out as coordinate vectors with respect to the standard basis. We already know that if we have a vector in the standard basis then $T(\vec{x}) = A\vec{x}$. Thus, we need only remember what a coordinate vector in the basis $\mathcal{B}$ means:

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}_\mathcal{B} = (\vec{b}_1)_\mathcal{S} \cdot x_1 + (\vec{b}_2)_\mathcal{S} \cdot x_2 + \cdots + (\vec{b}_m)_\mathcal{S} \cdot x_m$$

Thus, we have that $\vec{x}_\mathcal{S} = [\vec{b}_1 \cdots \vec{b}_m]\vec{x}_\mathcal{B}$. Hence we need only use the matrix $B = [\vec{b}_1 \cdots \vec{b}_m]$ to change to the standard basis, where we know which matrix represents $T$, and then $B^{-1}$ to change back to the basis $\mathcal{B}$. This gives:

$$T(\vec{x}_\mathcal{B}) = B^{-1}AB\vec{x}_\mathcal{B}$$

Thus, given any two bases, $\mathcal{B}_1 = \{\vec{b}_1, \vec{b}_2, ..., \vec{b}_m\}$ and $\mathcal{B}_2 = \{\vec{v}_1, \vec{v}_2, ..., \vec{v}_m\}$, and the change of basis matrix, $B$ such that

$$\vec{x}_{\mathcal{B}_1} = B\vec{x}_{\mathcal{B}_2},$$

if $T$ acts on coordinated vectors with respect to the basis $\mathcal{B}_1$ by the matrix $A$, then $T$ acts on coordinate vectors with respect to the basis $\mathcal{B}_2$ by the matrix $B^{-1}AB$.

**Question 1.3:** Might there be a BEST basis to represent a linear transformation by? If we want to analyze how a linear transformation acts on a space, is there a basis that makes it EASY see what is going on?

**Answer:** Yes. There is a "best basis" on which $T$ acts simply. Let's explore this in the next section.

## 1.1  Eigenvalues and Eigenvectors

**Let's make life easy and only consider $T : \mathbb{R}^n \to \mathbb{R}^n$ This way, we can view $T$ as a transformation of $\mathbb{R}^n$, rather than as a map between spaces. Since $T$ is a linear transformation, we know that it is represented by the action of some matrix, $A$.**

$$T(\vec{x}) = A\vec{x}$$

**Question 1.4:** If we want to find a basis for $\mathbb{R}^n$ on which $T$ acts simply, what might "simply" mean?

**Answer:** By the definition of vector spaces, there are only two things we can do with vectors:

1. We can scale vectors and get another vector.

2. We can add vectors together and get another vector.

The "simplest" is probably just scaling. Therefore, we are looking for vectors, $\vec{x} \in \mathbb{R}^n$ such that there is some scalar, $\lambda$, for which

$$T(\vec{x}) = \lambda\vec{x}.$$

This gives rise to the following definitions.

**Definition 1.2.** An **eigenvalue** of a matrix, $A$, is scalar, $\lambda$, such that there exists a non-zero vector for which

$$A\vec{x} = \lambda\vec{x}$$

Such a non-zero vector, $\vec{x}$, is called an **eigenvector**.

**Note:** Eigenvectors are associated to a particular eigenvalue. A vector cannot be an eigenvector for two distinct eigenvalues. Further, any scalar multiple of an eigenvector is also an eigenvector, since

$$A(c\vec{x}) = c(A\vec{x}) = c(\lambda\vec{x}) = \lambda(c\vec{x}).$$

We can rearrange the definition of an eigenvalue to read as follows.

**Definition 1.3.** An **eigenvalue** of a matrix, $A$, is scalar such that $null(A - \lambda I)$ is non-trivial. That is, there exist non-zero vectors, $\vec{x}$, called **eigenvectors** of A with eigenvalue $\lambda$, such that

$$(A - \lambda I)\vec{x} = \vec{0}.$$

**Definition 1.4.** The **eigenspace** associated to an eigenvalue, $\lambda$, of a matrix, $A$ is the space of all eigenvectors of that eigenvalue, plus the zero vector. That is, it is the set $null(A - \lambda I)$.

Observe that this means that the eigenspace of an eigenvalue is a vector space.

**Question 1.5:** Ok, now that we know that the vectors we want are called "eigenvectors," how do we find them??

**Answer:** First, we find eigenvalues, and then we find the eigenvectors. Since eigenvalues are scalars, $\lambda$, for which the matrix $(A - \lambda I)$ is NOT invertible, we know that **if $\lambda$ is an eigenvalue, then**

$$det(A - \lambda I) = 0$$

**Example 1.1:** Let $A = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}$. Then,

$$
\begin{aligned}
det(A - \lambda I) &= det \begin{pmatrix} 1 - \lambda & 2 \\ 0 & 1 - \lambda \end{pmatrix} \\
&= (1 - \lambda)^2
\end{aligned}
$$

Thus, the $det(A - \lambda I) = 0$ when $\lambda = 1$. $\lambda = 1$ is our only eigenvalue. To find eigenvectors associated to $\lambda = 1$, we need to find non-zero vectors which satisfy

$$(A - (1)I)\vec{x} = \begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix} \vec{x} = 0.$$

By 308 techniques, $\vec{x} = c \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, is a solution for any non-zero scalar, $c$. Thus, the vector $\vec{x} = c \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, is an eigenvector associated to the eigenvalue $\lambda = 1$ for any non-zero scalar, $c$.

**Definition 1.5.** The expression $det(A - \lambda I)$ is called the **characteristic polynomial**. The equation $det(A - \lambda I) = 0$ is called the **characteristic equation**.

**Question 1.6:** If we want to find all the eigenvalues, how many eigenvalues can an $(n \times n)-$matrix have?

**Answer:** If $A$ is an $(n \times n)-$matrix, then the characteristic polynomial, $det(A - \lambda I)$, is a polynomial in $\lambda$ of degree $n$. Thus, there are $n$ roots. That is, **counting multiplicity** there are $n$ solutions to the characteristic equation.

In Example 1.1, $\lambda = 1$ was the only eigenvalue, but it occurred as a double root of the the characteristic polynomial. Because of this, we count it twice. We make this rigorous by the following definition.

**Definition 1.6.** Given the characteristic polynomial, $det(A - \lambda I)$, we can factor it into a product of monomials,

$$det(A - \lambda I) = (\lambda - \lambda_1)^{m_1} \cdot (\lambda - \lambda_2)^{m_2} \cdots (\lambda - \lambda_j)^{m_j}$$

where each of the $\lambda_1, \lambda_2, ...., \lambda_j$ are distinct.

The **algebraic multiplicity** of an eigenvalue, $\lambda_i$, is the power, $m_i$, which shows up in this factorization of the characteristic polynomial. It is the "number of times" $\lambda_i$ is a root of the characteristic polynomial. We will often denote the algebraic multiplicity of an eigenvalue, $\lambda$, of $A$ by $m_\lambda$.

**Observation:** For an $(n \times n)-$matrix, $A$, if $m_1, m_2, ..., m_j$ are the algebraic multiplicities of the distinct eigenvalues of $A$, then

$$\sum_{i=1}^{j} m_i = n.$$

**Question 1.7:** If we want to find eigenvectors associated to an eigenvalue, how large can the eigenspace of an eigenvalue be?

**Definition 1.7.** If $A$ is an $(n \times n)-$matrix and $\lambda$ is an eigenvalue of A, then we know that $(A - \lambda I)$ is not invertible. The nullity of $(A - \lambda I)$, written $nullity(A - \lambda I)$, is the **geometric multiplicity** of $\lambda$. For a given eigenvalue, $\lambda$, of a matrix $A$ we will often denote the geometric multiplicity by $q_\lambda$.

**Answer:** It is a fact that for every eigenvalue, $\lambda$,

$$1 \leq q_\lambda \leq m_\lambda.$$

That is, the dimension of the eigenspace of an eigenvalue can be anywhere between 1 and the algebraic multiplicity of that eigenvalue.

The geometric multiplicity tells us how many linearly independent eigenvectors are associated to a given eigenvalue. Since we want to find a basis of eigenvectors, this number is very important. It gives rise to the following definitions.

Eigenvalues and their multiplicities tell us a lot about a matrix, and hence a lot about linear transformations. Because they are so important, we classify matrices according to their multiplicities.

**Definition 1.8.** A square matrix, $A$, with the property that for **every** eigenvalue $q_\lambda = m_\lambda$ is called **diagonalizable**.

**Definition 1.9.** A square matrix, $A$, for which there is an eigenvalue for which $q_\lambda < m_\lambda$ is called **defective**.

## 1.2 Jordan Normal Form

**So what was the point?** Remember, we are trying to find a basis for $\mathbb{R}^n$ such that our linear transformation $T : \mathbb{R}^n \to \mathbb{R}^n$, acts by scaling, because scaling is simple and easy to understand. In such a basis (if it can be found) the matrix which represents $T$ should be very simple and easy to analyze.

**Let us suppose that an $(n \times n)-$matrix, $A$, is diagonalizable.** Then,

$$\sum_{i=1}^{j} q_i = \sum_{i=1}^{j} m_i = n.$$

Thus, we can find $n$ linearly independent vectors, $\{\vec{\xi_1}, \vec{\xi_2}, ..., \vec{\xi_n}\}$, where each vector is an eigenvector of $A$. I claim that this is the "best" basis in which to view the linear transformation represented by $A$ Let's see why in an example.

**Example 1.2.** Let $A = \begin{pmatrix} 2 & -2 & 1 \\ -1 & 3 & -1 \\ 2 & -4 & 3 \end{pmatrix}$ represent a linear transformation,

$$T : \mathbb{R}^3 \to \mathbb{R}^3,$$

in the standard basis. That is, for all vectors $\vec{x} \in \mathbb{R}^3$, if we write $\vec{x}$ as a coordinate vector with respect to the standard basis,

$$\vec{x} = x_1\vec{e_1} + x_2\vec{e_2} + x_3\vec{e_3} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}_S,$$

then $T(\vec{x}) = \begin{pmatrix} 2 & -2 & 1 \\ -1 & 3 & -1 \\ 2 & -4 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}_S$ .

Let's find a better basis, a basis of eigenvectors of $A$, and see what $T$ is represented by in that basis!

To do this, we need to find eigenvalues and eigenvectors. We begin with the characteristic polynomial, $det(A - \lambda I) = (\lambda - 1)^2(\lambda - 6)$.

To find eigenvectors, we solve $(A - \lambda I)\vec{x} = 0$. For $\lambda = 1$, we get $\vec{\xi_1} = \begin{pmatrix} 2 \\ 1 \\ 0 \end{pmatrix}$ and

$\vec{\xi_2} = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$ . For $\lambda = 6$, we get $\vec{\xi_3} = \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}$ .

It is easy to check that $\{\vec{\xi_1}, \vec{\xi_2}, \vec{\xi_3}\}$ are linearly independent. This is left to the reader.

So, I claim that $T$ is represented in the basis $\mathcal{B} = \{\vec{\xi_1}, \vec{\xi_2}, \vec{\xi_3}\}$ by a simple matrix. Let $J$ be the matrix which represents $T$ in the basis $\mathcal{B}$.

**What is this matrix $J$?** Let's check by calculating $J$ directly! In the basis, $\mathcal{B} = \{\vec{\xi_1}, \vec{\xi_2}, \vec{\xi_3}\}$, the linear transformation $T$ is represented by some matrix, $J$. That means that if we write a vector, $\vec{x} \in \mathbb{R}^3$ as a coordinate vector in the basis, $\mathcal{B}$,

$$\vec{x} = x_1\vec{\xi_1} + x_2\vec{\xi_2} + x_3\vec{\xi_3} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}_\mathcal{B},$$

then we have that

$$
\begin{aligned}
T(\vec{x}) &= T(x_1\vec{\xi_1} + x_2\vec{\xi_2} + x_3\vec{\xi_3}) \\
&= x_1 T(\vec{\xi_1}) + x_2 T(\vec{\xi_2}) + x_3 T(\vec{\xi_3}) \quad &(1) \\
&= x_1 A\vec{\xi_1} + x_2 A\vec{\xi_2} + x_3 A\vec{\xi_3} \quad &(2) \\
&= x_1 (1)\vec{\xi_1} + x_2 (1)\vec{\xi_2} + x_3 (6)\vec{\xi_3} \quad &(3) \\
&= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 6 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}_{\mathcal{B}}. \quad &(4)
\end{aligned}
$$

Thus, $J = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 6 \end{pmatrix}$. Which is MUCH simpler than $A$, as promised.

Or, we could have remembered our answer to Question 1.2. That is, we could just do a change of basis transformation:

$$
J = \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\xi_3} \end{pmatrix}^{-1} A \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\xi_3} \end{pmatrix}
$$

.

For ease of computation, we rearrange the equation:

$$
\begin{aligned}
\begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\xi_3} \end{pmatrix} J &= A \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\xi_3} \end{pmatrix} \\
&= \begin{pmatrix} (1)\vec{\xi_1} & (1)\vec{\xi_2} & (6)\vec{\xi_3} \end{pmatrix} \\
&= \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\xi_3} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 6 \end{pmatrix}.
\end{aligned}
$$

Thus, if $A$ is diagonalizable, we can do a change of basis transformation,

$$
A = SJS^{-1}
$$

where $S$ is a matrix whose columns are a basis of eigenvectors and $J$ is a diagonal matrix with the eigenvalues on the diagonal.

**Note:** Note that $J$ acts by scaling the components. This makes sense because $T$ acts be scaling the eigenvectors. Also, the form of $J$ is very special. It is block diagonal, and **the eigenvalues occur on the diagonal as many times as their algebraic multiplicities.** This will come up again, later.

**Question 1.8:** What if $A$ is defective?

**Answer:** If $A$ is defective, then there exists at least one eigenvalue $\lambda$ for which $q_\lambda < m_\lambda$. Hence,

$$\sum_{i=1}^{j} q_i < \sum_{i=1}^{j} m_i = n.$$

Consider our earlier example $A = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}$. We could only find one linearly independent eigenvector.

Thus, $T$ is a linear transformation represented in some basis by the action of $A$, **if $A$ is defective, then we cannot find a basis of eigenvectors.** Observe that if we cannot find a basis of eigenvectors, then we cannot find a basis in which the linear transformation, $T$, is represented by a diagonal matrix. But, if we cannot represent $T$ by a diagonal matrix, we can represent it by an *almost* diagonal matrix.

The idea is very much the same: we want to choose a basis in which the action of $T$ is simple. Eigenvectors, which $T$ just rescales, are still the best. But if $A$ is not diagonalizable, we cannot find enough of them. Thus, we need some extra vectors on which $T$ acts in a simple way to complete our basis.

**Question 1.9:** What is another simple way in which $T$ can act on a vector?

**Answer:** The other "simple" way to act on a vector was to add some other vector to it. So, if we collect as many linearly independent eigenvectors as we can find, say, $\{\vec{\xi}_1, \vec{\xi}_2, , \vec{\xi}_m\}$, and we cannot find another vector, $\vec{\eta}$, such that $A\vec{\eta} = \lambda\vec{\eta}$ AND the collection, $\{\vec{\xi}_1, \vec{\xi}_2, , \vec{\xi}_m, \vec{\eta}\}$ is linearly independent, perhaps we could find a vector, $\vec{\eta}$, such that

$$A\vec{\eta} = \lambda\vec{\eta} + \vec{x}$$

and the collection, $\{\vec{\xi}_1, \vec{\xi}_2, , \vec{\xi}_m, \vec{\eta}\}$ is linearly independent.

**Question 1.10:** What $\vec{x}$ should we choose for $A\vec{\eta} = \lambda\vec{\eta} + \vec{x}$?

**Answer:** We could choose any $\vec{x}$, but we since we already have eigenvectors, and we know that $T$ acts simply on them, it is natural to see what happens if we try to find $\vec{\eta}$ such that

$$A\vec{\eta} = \lambda\vec{\eta} + \vec{\xi}$$

for $\vec{\xi}$ an eigenvector associated to the eigenvalue, $\lambda$.

Notice that we can rearrange this equation to read,

$$(A - \lambda I)\vec{\eta} = \vec{\xi}.$$

16

If we multiply both sides of the equation by the matrix, $(A - \lambda I)$, we see that $\vec{\eta}$ satisfies,

$$(A - \lambda I)^2 \vec{\eta} = (A - \lambda I)\vec{\xi} = \vec{0}$$

This gives rise to the following definitions.

**Definition 1.10.** A **generalized eigenvector of order m** of an eigenvalue, $\lambda$, of a matrix, $A$, is a non-zero vector, $\vec{\eta}$, such that

$$(A - \lambda I)^m \vec{\eta} = \vec{0}$$

but

$$(A - \lambda I)^{m-1} \vec{\eta} \neq \vec{0}$$

In practice, if might be very difficult to compute $(A - \lambda I)^m$ and $(A - \lambda I)^{m-1}$.

**Definition 1.11.** A **Jordan Chain of length m** associated to an eigenvalue, $\lambda$, of the matrix $A$ is an ordered collection of $m$ generalized eigenvectors, $\{\vec{\xi}, \vec{\eta}, ...., \vec{\rho}, \vec{\nu}\}$ such that

$$
\begin{aligned}
(A - \lambda I)\vec{\nu} &= \vec{\rho} \\
&\vdots \\
(A - \lambda I)\vec{\eta} &= \vec{\xi} \\
(A - \lambda I)\vec{\xi} &= \vec{0}
\end{aligned}
$$

From the definitions, it is clear that $\vec{\xi}$ is an eigenvector, $\vec{\eta}$ is a generalized eigenvector of order 2, $\vec{\rho}$ is a generalized eigenvector of order $m - 1$, and $\vec{\nu}$ is a generalized eigenvector of order m.

**The "best" basis that we are looking for will consist of eigenvectors and generalized eigenvectors arranged in a sequence of Jordan chains.**

We saw above that for diagonalizable matrices, we can do a change of basis and represent them by a diagonal matrix. Let's see what happens for defective matrices.

**Example 1.3** Let's return to our previous example, a linear transformation represented in the standard basis by the matrix $A = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}$. As we already saw, $A$ has one eigenvalue, $\lambda = 1$, with algebraic multiplicity 2. Let our eigenvector associated to $\lambda = 1$ be $\vec{\xi} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$.

To complete the basis, we need to build a Jordan chain of length 2 off of $\vec{\xi}$. That is, we need to solve $(A - (1)I)\vec{\eta} = \vec{\xi}$.

$$\begin{pmatrix} 0 & 2 \\ 0 & 0 \end{pmatrix} \vec{\eta} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Thus, we can choose $\vec{\eta} = \begin{pmatrix} 0 \\ 1/2 \end{pmatrix}$. Thus, our basis is $\mathcal{B} = \{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1/2 \end{pmatrix} \}$

By change of basis, then, $T$ is represented in the basis $\mathcal{B}$ by the action of the matrix $\left( \vec{\xi} \ \ \vec{\eta} \right)^{-1} A \left( \vec{\xi} \ \ \vec{\eta} \right)$. Explicitly calculating this matrix, we get,

$$
\begin{aligned}
\left( \vec{\xi} \ \ \vec{\eta} \right) J &= A \left( \vec{\xi} \ \ \vec{\eta} \right) \\
&= \left( A\vec{\xi} \ \ A\vec{\eta} \right) \\
&= \left( (1)\vec{\xi} \ \ (1)\vec{\eta} + \vec{\xi} \right) \\
&= \left( \vec{\xi} \ \ \vec{\eta} \right) \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}
\end{aligned}
$$

Hence, $J = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$. **Note that even though $J$ is not diagonal, it still is very simple. It has the eigenvalues on the diagonal and a 1 on the super-diagonal.** Also, it acts simply. In the first coordinate, it simply scales by the eigenvalue 1.

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ 0 \end{pmatrix} = \begin{pmatrix} (1)x_1 \\ 0 \end{pmatrix}$$

In the second coordinate, it scales the original vector by the eigenvalue 1 and "rotates" the vector by adding something in the first coordinate.

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ x_2 \end{pmatrix} = \begin{pmatrix} (1)x_2 \\ (1)x_2 \end{pmatrix}$$

This completely describes the action of the linear transformation $T$. We generalize this form of matrix with the following defintion.

**Definition 1.12.** A square matrix, $J$, is called a **Jordan Block** if $J$ has the following properties:

1. $J$ has the same value in all the diagonal entries.

2. All of the entries on the super-diagonal are 1.

3. All the other entries are 0

**Examples.** The following are Jordan Blocks: $\begin{pmatrix} 5 & 1 \\ 0 & 5 \end{pmatrix}$, $\begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$, $(13)$.

The following are NOT Jordan Blocks:
$$\begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}.$$

**Definition 1.13.** Let $A$ be an $(n \times n)-$ matrix. The **Jordan Form of** $A$ is a block-diagonal $(n \times n)-$matrix with the following properties:

1. Each block on the diagonal is a Jordan block with an eigenvalue of $A$ on the diagonal.

2. Each eigenvalue, $\lambda$, has $q_\lambda$ Jordan blocks, arranged adjacently.

3. For each eigenvalue, $\lambda$, the sum of the sizes of the Jordan blocks associated to that eigenvalue adds up to $m_\lambda$.

4. There is a change of basis matrix, $S$, such that
$$AS = SJ$$

.

**Examples.** The following are examples of matrices in Jordan Normal Jorm:
$$\begin{pmatrix} 4 & 0 \\ 0 & 5 \end{pmatrix}, \begin{pmatrix} 3 & 1 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & -17 \end{pmatrix}, \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -17 & 0 \\ 0 & 0 & 0 & -17 \end{pmatrix}.$$

The following are NOT in Jordan Normal Form:
$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \begin{pmatrix} 3 & 1 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 1 & -17 \end{pmatrix}.$$

**The Jordan form of a matrix is the "simple" matrix we are looking for.**

**Question 1.11:** So where do Jordan chains come in?

**Answer:** Each Jordan block in $J$ corresponds to a Jordan chain in $S$. Consider the following example. Let the matrix, $A$, have the Jordan form,
$$J = \begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix}$$

19

Let $\vec{\xi}$ be an eigenvector of $A$ associated to $\lambda$. We know that there is a change of basis matrix, $S$, such that $AS = SJ$. As mentioned, the columns of $S$ will be the "best" basis in which to represent $T$. We already know we want $\vec{\xi}$ to be a vector in $S$, so let us write

$$S = \begin{pmatrix} \vec{\xi} & \vec{\eta} \end{pmatrix}$$

and see what properties $\vec{\eta}$ must have.

On the one hand,

$$\begin{aligned} AS &= A \begin{pmatrix} \vec{\xi} & \vec{\eta} \end{pmatrix} \\ &= \begin{pmatrix} A\vec{\xi} & A\vec{\eta} \end{pmatrix} \\ &= \begin{pmatrix} \lambda\vec{\xi} & A\vec{\eta} \end{pmatrix} \end{aligned}$$

And on the other hand,

$$\begin{aligned} SJ &= \begin{pmatrix} \vec{\xi} & \vec{\eta} \end{pmatrix} \begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix} \\ &= \begin{pmatrix} \lambda\vec{\xi} & \lambda\vec{\eta} + \vec{\xi} \end{pmatrix} \end{aligned}$$

Thus, it must be that $A\vec{\eta} = \lambda\vec{\eta} + \vec{\xi}$. Rearranging, we have that $(A - \lambda I)\vec{\eta} = \vec{\xi}$. Thus, a $(2 \times 2)$−Jordan block in $J$ corresponds to a Jordan chain of length 2 in $S$. The same goes for larger blocks and chains. A $(3 \times 3)$−Jordan block corresponds to a Jordan chain of length 3, etc.

Putting it all together, if the Jordan form of a matrix, $A$, has a $(1 \times 1)$−Jordan block, a $(3 \times 3)$−Jordan block, and a $(2 \times 2)$−Jordan block down the diagonal, then $S$ has– from left to right– a Jordan chains of length 1, Jordan chains of length 3, and Jordan chains of length 2.

Looking at the definition of the Jordan form, we see that **for $S$, each linearly independent eigenvector we found will be the base of a Jordan chain, and the sum of the lengths Jordan chains associated to a given eigenvalue, $\lambda$, must add up to $m_\lambda$, the algebraic multiplicity of that eigenvalue. Therefore, we only need find generalized eigenvectors and build Jordan chains for eigenvalues, $\lambda$, for which $q_\lambda < m_\lambda$.**

**Question 1.12** Suppose that the Jordan form of a matrix, $A$, has all $(1 \times 1)$−Jordan blocks on the diagonal. Is $A$ diagonalizable or defective?

**Answer:** If $J$ has only $(1 \times 1)$−Jordan blocks, then $S$ has only Jordan chains of length 1. Since the first vector on every Jordan chain is an eigenvector, $S$ is only eigenvectors. Therefore, $A$ must have $n$ linearly independent eigenvectors, and hence $A$ is diagonalizable. This should make sense because if $J$ has only $(1 \times 1)$−Jordan blocks, then $J$ is a diagonal matrix!

## 1.3  How to find the Jordan Form of a matrix $A$

In this class we will only do computations with $(2 \times 2)$ and $(3 \times 3)$ matrices. This actually simplifies our work greatly. In fact, to find the Jordan Form of a matrix which is $(3 \times 3)$ or smaller, we need only look at the algebraic and geometric multiplicities. To find the change of basis matrix, $S$, however, we will need to do a bit of calculation.

**Example 1.4** Let $A$ be a matrix with eigenvalues $\lambda = 2$ and $\lambda = 5$. Let $m_2 = 1$, $q_2 = 1$, $m_5 = 2$, and $q_5 = 1$.

Since we know that the algebraic multiplicities add up to the size of the matrix, we see that $A$ must be a $(3 \times 3)-$matrix. Thus, it's Jordan Form, $J$, must also be a $(3 \times 3)-$ matrix. This Jordan form will have 2 Jordan blocks on the diagonal, because there are only 2 linearly independent eigenvectors– one associated to $\lambda = 2$ and one associated to $\lambda = 5$. The sizes of the blocks associated to a distinct eigenvector must sum up to the algebraic multiplicity of that eigenvector, so the Jordan block associated to $\lambda = 2$ is a $(1)-$block. The Jordan block associated to $\lambda = 5$ is a $(2)-$block. Thus, the Jordan Form of $A$ must be either

$$
J = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 5 & 1 \\ 0 & 0 & 5 \end{pmatrix} \quad \text{or} \quad J = \begin{pmatrix} 5 & 1 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix}.
$$

Either is correct. **The Jordan Form of a matrix is not unique.**

**Example 1.5** Let $A$ be a matrix with eigenvalues $\lambda_1 = 7$ and $\lambda_2 = 0$. Let $m_7 = 1$, $q_7 = 1$, $m_0 = 2$, and $q_0 = 2$.

By the same reasoning as before, $A$ must again be a $(3 \times 3)-$matrix. Thus, it's Jordan Form, $J$, must also be a $(3 \times 3)-$ matrix. This Jordan form will have 3 Jordan blocks on the diagonal, because the sum of the geometric multiplicities is 3. That is, we have 3 linearly independent eigenvectors– one associated to $\lambda = 7$ and two associated to $\lambda = 0$. The sizes of the blocks associated to a distinct eigenvector must sum up to the algebraic multiplicity of that eigenvector, so the Jordan block associated to $\lambda = 7$ is a $(1 \times 1)-$block. The sizes of the two Jordan blocks associated to $\lambda = 0$ must add up to 2. But there is only one way to do this, with two $(1 \times 1)-$blocks. Thus, the Jordan Form of $A$ must be either

$$
J = \begin{pmatrix} 7 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{or} \quad J = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 7 \end{pmatrix}.
$$

Again, either is correct.

If we want to find $S$, however, we need to know the entries of $A$ and do actual calculations.

**Example 1.6** Let $T : \mathbb{R}^3 \to \mathbb{R}^3$ be a linear transformation represented by the matrix

21

$A = \begin{pmatrix} -3 & 3 & 6 \\ 0 & 1 & 2 \\ 0 & 4 & -1 \end{pmatrix}$ . Find the, $J$, Jordan Form of $A$ and a basis for $\mathbb{R}^3$ in which $T(\vec{x}) = J\vec{x}$.

- **Step 1: Find the eigenvalues and their algebraic multiplicities.** To find eigenvalues, we factor the characteristic polynomial, $det(A - \lambda I)$.

$$\begin{aligned} \det(A - \lambda I) &= (-3 - \lambda)[(1 - \lambda)(-1 - \lambda) - (2)(4)] \\ &= -(\lambda + 3)[\lambda^2 - 1 - 8] \\ &= -(\lambda + 3)^2(\lambda - 3) \end{aligned}$$

The eigenvalues are $\lambda = 3$ and $\lambda = -3$, with algebraic multiplicities 1 and 2, respectively.

- **Step 2: Find the eigenvectors.** Now we must find non-trivial solutions to $(A - \lambda I)\vec{x} = \vec{0}$. By 308, this means we need to row-reduce $(A - \lambda I)$.

For $\lambda = 3$, we calculate as follows

$$\begin{aligned} A - (3)I &= \begin{pmatrix} -6 & 3 & 6 \\ 0 & -2 & 2 \\ 0 & 4 & -4 \end{pmatrix} \\ &\rightarrow \begin{pmatrix} -6 & 0 & 9 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \end{aligned}$$

Since this matrix has one free variable, the geometric multiplicity of $\lambda = 3$ is 1. Choose the eigenvector $\vec{\xi_1} = \begin{pmatrix} 3 \\ 2 \\ 2 \end{pmatrix}$

For $\lambda = -3$, we calculate

$$\begin{aligned} A - (-3)I &= \begin{pmatrix} 0 & 3 & 6 \\ 0 & 4 & 2 \\ 0 & 4 & 2 \end{pmatrix} \\ &\rightarrow \begin{pmatrix} 0 & 2 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \end{aligned}$$

Since this matrix has one free variable, the geometric multiplicity of $\lambda = -3$ is 1. Choose the eigenvector $\vec{\xi_2} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$.

- **Step 3: Write down the Jordan Form.** The previous calculation shows that both eigenvalues have geometric multiplicity 1. Since the sum of the geometric multiplicities $(1 + 1 = 2)$ is equal to the number of Jordan blocks in $J$ (since each eigenvector is the

base of a Jordan chain) we know that there are two Jordan blocks in $J$. Since $\lambda = -3$ has algebraic multiplicity 2, it must show up on the diagonal twice. Thus, we have a $(1 \times 1)-$Jordan block associated to $\lambda = 3$ and a $(2 \times 2)-$Jordan block associated to $\lambda = -3$. Choose

$$J = \begin{pmatrix} 3 & 0 & 0 \\ 0 & -3 & 1 \\ 0 & 0 & -3 \end{pmatrix}.$$

- **Step 4: Build Jordan chains, if necessary.** Since our matrix $A$ is defective, we need to complete $\{\vec{\xi_1}, \vec{\xi_2}\}$ into a basis. This means building a Jordan chain off $\vec{\xi_2}$, because that one is associated to the eigenvalue with $q_\lambda < m_\lambda$.

  Hence, we solve $(A - (-3)I)\vec{\eta} = \vec{\xi_2}$,

$$\begin{pmatrix} 0 & 3 & 6 \\ 0 & 4 & 2 \\ 0 & 4 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

  Row-reducing, we get that

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = s \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ -1/9 \\ 2/9 \end{pmatrix}$$

  Thus, we choose our generalized eigenvector to be $\vec{\eta} = \begin{pmatrix} 0 \\ -1/9 \\ 2/9 \end{pmatrix}$.

  This gives us two Jordan chains, $\{\vec{\xi_1}\}$ and $\{\vec{\xi_2}, \vec{\eta}\}$.

- **Step 5: Put your Jordan chains in the right order.** We chose $J = \begin{pmatrix} 3 & 0 & 0 \\ 0 & -3 & 1 \\ 0 & 0 & -3 \end{pmatrix}$.

  Going down the diagonal, we have a $(1 \times 1)-$block and then a $(2 \times 2)-$block. Thus, to make the equation $AS = SJ$ hold, from left to right, the columns of $S$ must form a Jordan chain of length 1 and then a Jordan chain of length 2. That is,

$$\begin{aligned} A \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\eta} \end{pmatrix} &= \begin{pmatrix} (3)\vec{\xi_1} & (-3)\vec{\xi_2} & (-3)\vec{\eta} + \vec{\xi_2} \end{pmatrix} \\ &= \begin{pmatrix} \vec{\xi_1} & \vec{\xi_2} & \vec{\eta} \end{pmatrix} \begin{pmatrix} 3 & 0 & 0 \\ 0 & -3 & 1 \\ 0 & 0 & -3 \end{pmatrix}. \end{aligned}$$

With a few caveats, this process generalizes to all $(2 \times 2)-$ and $(3 \times 3)-$matrices. **Steps 1 and 2 are always the same for all matrices.** As noted earlier, the Jordan Form is not unique, so there are some choices involved in which Jordan Form you write down for Step 3.

**The Caveats**

1. If $A$ is diagonal, then we already found a basis in Step 2, so Step 4 is unneccessary, since there are no generalized eigenvectors to be found.

2. If $A$ is defective, then things can get a little tricky. **We need the lengths of the Jordan chains of each eigenvalue to sum up to the algebraic multiplicity of that eigenvalue.** Thus, we need only build Jordan chains for eigenvectors, $\vec{\xi}$, associated to eigenvalues, $\lambda$, with $q_\lambda < m_\lambda$.

Case 1. If $A$ is a $(2 \times 2)-$matrix, then this is easy. We already have one eigenvector, $\vec{\xi}$, and **we can always solve** $(A - \lambda I)\vec{\eta} = \vec{\xi}$ for $\vec{\eta}$.

Case 2. If $A$ is a $(3 \times 3)-$ matrix, then there are only a few possibilities. If $A$ has two distinct eigenvalues, then just as in the example above, it must be that they have $m_1 = 1$, $m_2 = 2$, $q_1 = 1$, and $q_2 = 1$. **We can always solve** $(A - \lambda I)\vec{\eta} = \vec{\xi_2}$ for $\vec{\eta}$ as we did, above.

Case 3. If $A$ has one eigenvalue with $m = 3$ and $q = 1$, then **we can always build our Jordan chain by solving**

$$(A - \lambda I)\vec{\xi} = 0$$
$$(A - \lambda I)\vec{\eta} = \vec{\xi}$$
$$(A - \lambda I)\vec{\nu} = \vec{\eta}$$

for $\vec{\eta}$ and then $\vec{\nu}$.

Case 4. If $A$ has only one eigenvalue with $m = 3$ and $q = 2$, however, then things get tricky. There will only be one eigenvector off of which you can build a Jordan chain. Further, **it is possible to make poor choices of both $\vec{\xi_1}$ and $\vec{\xi_2}$ so that it will not be possible to solve** $(A - \lambda I)\vec{\eta} = \vec{\xi}$ **for** $\vec{\eta}$. Consider the following example:

$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}$ and $\vec{\xi_1} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$ and $\vec{\xi_2} = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix}$. It is not possible to solve either

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} \vec{\eta} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} \vec{\eta} = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix}$$

In this example, one must choose $\vec{\xi} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$ or some scalar multiple of that eigenvector in order to build a Jordan chain.

In general, if you have chosen poorly, you many need to choose $\vec{\eta}$ first. That is, choose **any vector outside** $span\{\vec{\xi_1}, \vec{\xi_2}\}$, let this vector be $\vec{\eta}$ and then solve for a new eigenvector by solving $(A - \lambda I)\vec{\eta} = \vec{\xi_0}$ for $\vec{\xi_0}$. This gives a Jordan chain

of length 2. Choosing whichever vector, $\vec{\xi}_1$ or $\vec{\xi}_2$, is not a scalar multiple of your new $\vec{\xi}_0$ to be the Jordan chain of length 1 solve the problem. **In this class, this will be very, very rare. Usually, one of the obvious choices for $\vec{\xi}_1$ and $\vec{\xi}_2$ will work.**

3. For Step 5, once the Jordan Form you want to use has been chosen, the order of the Jordan chains in the change of basis matrix $S$ is fixed. Going left to right, then must correspond in both size and eigenvalue with the Jordan blocks on the diagonal of $J$, otherwise the equation $AS = SJ$ will not hold.

### Reflection Questions

1. In your own words, what information do you need in order to write down the Jordan form of a $(2 \times 2)$ or $(3 \times 3)$ matrix?

2. In your own words, what steps do we take to get the information needed in order to write down the the the Jordan form of a $(2 \times 2)$ or $(3 \times 3)$ matrix?

3. Let $A$ be a $(2 \times 2)$ or $(3 \times 3)$ matrix and $J$ the Jordan form of $A$. In your own words, what are the steps needed to write down a change of basis matrix, $S$, such that $AS = SJ$?

4. If $A$ is a matrix with two eigenvalues, $\lambda_1$, with $m_{\lambda_1} = 3$ and $q_{\lambda_1} = 3$ and $\lambda_2$, with $m_{\lambda_2} = 2$ and $q_{\lambda_2} = 1$. How big is $A$? How many Jordan blocks will the Jordan form of $A$ have? How many Jordan chains will $S$ have? How many columns will $S$ have?

5. If $AS = SJ$ for matrices $J = \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 1 \\ 0 & 0 & \lambda_2 \end{pmatrix}$ and $S = \begin{pmatrix} \vec{\xi}_1 & \vec{\xi}_2 & \vec{\nu} \end{pmatrix}$. If we choose a different Jordan form, $\tilde{J} = \begin{pmatrix} \lambda_2 & 1 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_1 \end{pmatrix}$, what must the corresponding $\tilde{S}$ be if $A\tilde{S} = \tilde{S}\tilde{J}$?

6. If $A$ is diagonalizable, will $J$ be a diagonal matrix?

7. If $\{\vec{\xi}, \vec{\nu}\}$ is a Jordan chain associated to an eigenvalue, $\lambda$, why must $\vec{\xi}$ and $\vec{\nu}$ be linearly independent?

# 2 Introduction to Linear Systems of Ordinary Differential Equations

In this class we will only be dealing with $1^{st}$ order linear ODEs. In 308, you spent a lot of time studying $2^{nd}$ order linear ODEs. While the techniques you learned are interesting, there is no need to go beyond $1^{st}$ order linear ODEs. **Every linear ODE of finite order can be reduced to a $1^{st}$ order linear system of ODEs.**

**Example 2.1.** Consider a general $2^{nd}$ order linear ODE: $y'' + b(t)y' + c(t)y = d(t)$. By allowing us to take advantage of the power of vectors and matrices, we can simply let the first coordinate represent $y$ and the second represent $y'$. This allows us to write the equation $y'' + b(t)y' + c(t)y = d(t)$ as

$$\begin{pmatrix} y \\ y' \end{pmatrix}' = \begin{pmatrix} y' \\ y'' \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -c(t) & -b(t) \end{pmatrix} \begin{pmatrix} y \\ y' \end{pmatrix} + \begin{pmatrix} 0 \\ d(t) \end{pmatrix}$$

which we may abbreviate as $\frac{d}{dt}\vec{y}(t) = A(t)\vec{y}(t) + \vec{g}(t)$

Recall that matrices and vectors are a book-keeping device to encode systems of equations. Thus, the above vector-valued differential equation can also be written as a system of first order ODES:

$$\begin{array}{rlll} y_1'(t) & = & & y_2(t) \\ y_2'(t) & = & -c(t)y_1(t) & -b(t)y_2(t) & +d(t) \end{array}$$

Note that the entries of $A(t)$ show up as the coefficients in the system of equations. In this class, we will only deal with matrices with constant coefficients. That is, we will only study constant-coefficient linear systems of ODEs.

## 2.1 Theorems and Definitions

**Question 2.1:** What does a constant coefficient, $1^{st}$ order linear system of ODEs look like??

**Answer:** In general, they look like

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + \vec{g}(t).$$

Where $\vec{x}(t)$ and $\vec{g}(t)$ are **vector-valued functions**,( i.e., vectors whose entries are scalar-valued functions ) and $A$ is a square matrix with constant entries.

**Definition 2.1.** A **solution to a constant-coefficient** $1^{st}$ **order linear system of ODEs,**

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + \vec{g}(t),$$

**on the interval** $(a, b)$ is a vector-valued function, $\vec{\phi}(t)$, such that for all $t$ in the interval $(a, b)$, $\vec{\phi}(t)$ satisfies the equation:

$$\frac{d}{dt}\vec{\phi}(t) = A\vec{\phi}(t) + \vec{g}(t).$$

In this class, we will usually only deal with **global solutions**, that is, solutions for all values of $t \in \mathbb{R}$.

**Definition 2.2.** A constant coefficient, $1^{st}$ order linear system of ODEs is called **homogeneous** if $\vec{g}(t) \equiv 0$. It is called **non-homogeneous** if $\vec{g}(t) \not\equiv 0$.

**Theorem 2.1.** *(Existence and Uniqueness) Let $A$ be a constant coefficient $(n \times n)$-matrix. Let $\vec{g}(t)$ be a continuous vector-valued function. For every $\vec{x}_0 \in \mathbb{R}^n$ and every $t_0 \in \mathbb{R}$, there exists a unique solution to the Initial Value Problem*

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + \vec{g}(t) \qquad \vec{x}(t_0) = \vec{x}_0.$$

This theorem is very important. It tells us that we have solutions and the solutions are unique. We will definitely use it later.

For now, let us restrict our attention to the homogeneous case.

**Question 2.2:** Now that we know we have solutions, what does the set of solutions to $\frac{d}{dt}\vec{x} = A\vec{x}$ look like? What kind of structure does it have?

**Answer:** To answer this question, it helps to re-write the equation $\frac{d}{dt}\vec{x} = A\vec{x}$ as the following:

$$(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}.$$

That is, we view solutions as vector-valued functions which are in the **kernel of the operator** $(\frac{d}{dt} - A)$.

**Definition 2.3.** An **operator** is a function of functions. That is, if is an object that takes functions as inputs and gives functions as outputs.

In our case, it is clear that for any differentiable vector-valued function, $\vec{y}(t)$, we can form a new function by the rule

$$(\frac{d}{dt} - A)\vec{y}(t) = \frac{d}{dt}\vec{y}(t) - A\vec{y}(t).$$

This leaves only the question of whether or not the operator $(\frac{d}{dt} - A)$ is linear. Let's check from the definition!

Let $\vec{x}(t)$ and $\vec{y}(t)$ be elements of the kernel of $(\frac{d}{dt} - A)$, then

$$
\begin{aligned}
(\tfrac{d}{dt} - A)(\vec{x}(t) + c\vec{y}(t)) &= (\tfrac{d}{dt} - A)\vec{x}(t) + c(\tfrac{d}{dt} - A)\vec{y}(t) \\
&= 0 + c0 \\
&= 0.
\end{aligned}
$$

**The kernel of a linear map is always a a vector space. Thus, the set of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$ forms a vector space.** This means that we can think about a solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$. as a vector in a vector space. That is, the set of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$— which is a collection of vector-valued functions— is a vector space. This is very useful, because in 308, we learned a lot about vector spaces.

**Question 2.3:** If solutions are vectors in a vector space, is there a notion of linear independence for solutions/vector-valued functions?

**Answer:** Yes. In fact, we have exactly the same definitions. We simply treat the vector-valued functions as different vectors at different times.

**Definition 2.4.** A collection of $n-$vector-valued functions, $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_j(t)\}$, is called **linearly dependent at time** $t$ if there exists a non-trivial linear combination such that

$$c_1\vec{x}_1(t) + c_2\vec{x}_1(t) + .... + c_j\vec{x}_j(t) = \vec{0}.$$

If the only such linear combination is the trivial combination (all zeros), the the collection is **linearly independent at time** $t$.

In practice, we will be mostly interested in when $j = n$ if $A$ is an $(n \times n)-$matrix. In that case, there is an easy test for linear independence.

**Definition 2.5.** For a collection of $n-$vector-valued functions, $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$, we define their **Wronskian** as the following function:

$$W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) = det \left( \vec{x}_1(t) \quad \vec{x}_2(t) \quad \cdots \quad \vec{x}_n(t) \right).$$

That is, the Wronskian of a collection of vector-valued functions is the determinant of the matrix whose columns are that collection of vector-valued functions. **This makes the Wronskian a function of** $t$. Just as for constant-vectors, if the $W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) = 0$ then the columns are linearly dependent at time $t$. If $W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) \neq 0$, then the columns are linearly independent at time $t$.

Now, at first, this should seem very unsatisfying. **It looks like vector-valued functions can be linearly independent at one time and linearly dependent at another time. In general, this does happen. But, not if** $\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)$ **are all solutions to** $(\frac{d}{dt} - A)\vec{x}(t) = 0$ **on an interval** $(a, b)$.

**Theorem 2.2.** *(Abel's Theorem) Let $A$ be an $(n \times n)-matrix$. Then, if $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$ are all solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$, then either*

$$W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) = 0 \qquad \textit{for all times } t$$

*or*

$$W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) \neq 0 \qquad \textit{for all times } t$$

*This means that solutions are either always linearly independent or always dependent.*

**Question 2.4:** If the space of solutions is a vector space, how big is it? What is its dimension?

**Answer: If $A$ is an $(n \times n)-$matrix, then the set of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$ is an $n-$dimensional vector space.**

Recall that one of the definitions of the dimension of a vector space is the number of vectors in a basis for that vector space. This is, it is the maximum number of linearly independent vectors that can be found in that vector space. **To show that we have an $n-$dimensional vector space, then, we must show that there exist $n$ linearly independent solutions, and that any collection of $n+1$ solutions is linearly dependent.**

To see that there exist at least $n$ solutions which are linearly independent, we turn to our Existence and Uniqueness Theorem. Choose time $t = 0$, for any vector, $\vec{x}_0 \in \mathbb{R}^n$, the Existence and Uniqueness Theorem says that there exists a unique solution to the Initial Value Problem :

$$(\frac{d}{dt} - A)\vec{x}(t) = 0 \quad , \quad \vec{x}(0) = \vec{x}_0.$$

So, we can choose any $n$ linearly independent vectors in $\mathbb{R}^n$, say, $\vec{x}_1, \vec{x}_2, ..., \vec{x}_n$ and the Existence and Uniqueness Theorem gives us $n$ solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$. For ease, let us call these solutions $\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)$. To see that these solutions are linearly independent, we check their Wronskian. By Abel's Theorem, we need only check at $t = 0$. Thus, we have that

$$
\begin{aligned}
W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](0) &= det \begin{pmatrix} \vec{x}_1(0) & \vec{x}_2(0) & \cdots & \vec{x}_n(0) \end{pmatrix} \\
&= det \begin{pmatrix} \vec{x}_1 & \vec{x}_2 & \cdots & \vec{x}_n \end{pmatrix} \\
&\neq 0
\end{aligned}
$$

Hence, we can find at least $n$ linearly independent solutions.

To see that any collection of $n + 1$ solutions must be linearly dependent, we also use the Existence and Uniqueness Theorem. Let $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t), \vec{x}_{n+1}(t)\}$ all be solutions to the equation $(\frac{d}{dt} - A)\vec{x}(t) = 0$. If the first $n$ of them are linearly dependent, then the whole collection is linearly dependent, and we have nothing to show. So, assume that $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$ are linearly independent. We need to show that $\vec{x}_{n+1}(t)$ is a linear combination of $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$.

Let $t = 0$. By Abel's Theorem, $\{\vec{x}_1(0), \vec{x}_2(0), ..., \vec{x}_n(0)\}$ is still a linearly independent collection of vectors in $\mathbb{R}^n$. Therefore, they must form a basis for $\mathbb{R}^n$. Therefore there is a linear combination such that

$$c_1\vec{x}_1(0) + c_2\vec{x}_2(0) + ... + c_n\vec{x}_n(0) = \vec{x}_{n+1}(0)$$

Since the solutions $\vec{x}_{n+1}(0)$ and $c_1\vec{x}_1(0) + c_2\vec{x}_2(0) + ... + c_n\vec{x}_n(0)$ satisfy the same Initial Value Problem, the Existence and Uniqueness Theorem implies that they must be the same solution. Thus,

$$\vec{x}_{n+1}(t) = c_1\vec{x}_1(t) + c_2\vec{x}_2(t) + ... + c_n\vec{x}_n(t).$$

This means that $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t), \vec{x}_{n+1}(t)\}$ is linearly dependent. Therefore, **the set of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$ is an $n-$dimensional vector space.**

**Question 2.5:** Now that we know that the space of solutions is an $n-$dimensional vector space, how do we get a handle on it? How do we describe it?

**Answer: All the information in a vector space is contained in its basis. Once we have a basis, we know everything about that vector space.** That means, once we have a basis for the space of solutions to the equation $(\frac{d}{dt} - A)\vec{x}(t) = 0$, we know everything about all solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$. Because bases are so important, they get a special name.

**Definition 2.6.** Let $A$ be an $(n \times n)-$matrix. A collection of vector-valued functions, $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$, is called a **fundamental set of solutions** for the system $(\frac{d}{dt} - A)\vec{x}(t) = 0$ if and only if the following hold:

1. Each vector-valued function $\vec{x}_i(t)$ is a solution of $(\frac{d}{dt} - A)\vec{x}(t) = 0$.

2. The collection, $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$, is linearly independent.

**A fundamental set of solutions of $(\frac{d}{dt} - A)\vec{x}(t) = 0$ is a basis for the space of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = 0$.**

Sometimes it is more convenient to put a fundamental set of solutions into a matrix. We have a special name for that, too.

**Definition 2.7.** Any $(n \times n)-$matrix-valued function $\Psi(t)$ whose columns form a fundamental set of solutions for $(\frac{d}{dt} - A)\vec{x}(t) = 0$ is called a **fundamental matrix** of the equation $(\frac{d}{dt} - A)\vec{x}(t) = 0$.

This means that the columns of a fundamental matrix, $\Psi(t)$, are solutions and are linearly independent. If we let $\{\vec{x}_1(t), \vec{x}_2(t), ..., \vec{x}_n(t)\}$ be a fundamental set of solutions for $(\frac{d}{dt} - A)\vec{x}(t) = 0$ and write $\Psi(t) = (\vec{x}_1(t) \quad \vec{x}_2(t) \quad \cdots \quad \vec{x}_n(t))$, then because the columns are solutions, we have that

$$
\begin{aligned}
A\Psi(t) &= A\left(\vec{x}_1(t) \quad \vec{x}_2(t) \quad \cdots \quad \vec{x}_n(t)\right) \\
&= \left(A\vec{x}_1(t) \quad A\vec{x}_2(t) \quad \cdots \quad A\vec{x}_n(t)\right) \\
&= \left(\frac{d}{dt}\vec{x}_1(t) \quad \frac{d}{dt}\vec{x}_2(t) \quad \cdots \quad \frac{d}{dt}\vec{x}_n(t)\right) \\
&= \frac{d}{dt}\Psi(t)
\end{aligned}
$$

This means that **any fundamental matrix to the equation $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$ is a solution to the corresponding matrix equation $(\frac{d}{dt} - A)Z(t) = 0_{n \times n}$.**

Because the columns of a fundamental matrix are linearly independent, for any fundamental matrix, $\Psi(t)$ , to the equation $(\frac{d}{dt} - A)\vec{x}(t) = 0$ we have that

$$
\begin{aligned}
det(\Psi(t)) &= det\left(\vec{x}_1(t) \quad \vec{x}_2(t) \quad \cdots \quad \vec{x}_n(t)\right) \\
&= W[\vec{x}_1, \vec{x}_2, ..., \vec{x}_n](t) \\
&\neq 0
\end{aligned}
$$

**Question 2.8:** Why do we care?

**Answer: Since the columns of a fundamental matrix form a basis for the set of solutions, every solution can be uniquely written as a linear combination of those basis vectors. Thus, if $\Psi(t)$ is a fundamental matrix for $(\frac{d}{dt} - A)\vec{x}(t) = 0$, then every solution to $(\frac{d}{dt} - A)\vec{x}(t) = 0$ can be written uniquely as**

$$\Psi(t)\vec{c}.$$

**for some vector $\vec{c} \in \mathbb{R}^n$.**

## 2.2   Matrix Exponentiation

Consider the scalar-valued ODE

$$y' = ay.$$

We know from 307 that the general solution to this ODE is

$$x(t) = e^{at} \cdot c.$$

Thus, we might hope that if we consider the vector-valued ODE

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}$$

the general solution might be something like

$$e^{At}\vec{c}$$

If we could make sense of it. That is, we might hope that $e^{At}$ was a fundamental matrix for $\frac{d}{dt}\vec{x} = A\vec{x}$.

**Question 2.6:** What could $e^{At}$ possibly mean??

**Answer:** To answer this question, we must ask what the scalar-valued function $e^{at}$ means. From 126, we recall that $e^{at}$ has a **Taylor series** or **power series**.

$$\begin{aligned}
e^{at} &= 1 + at + a^2\frac{t^2}{2!} + a^3\frac{t^3}{3!} + a^4\frac{t^4}{4!} + \dots \\
&= \sum_{n=1}^{\infty} a^n \frac{t^n}{n!}
\end{aligned}$$

This formal power series only involves taking powers of $a$, multiplying by scalars, $\frac{t^n}{n!}$, and summation. These are all things that we can do with matrices. So we arrive at the following definition.

**Definition 2.8.** Let $A$ be an $(n \times n)$−matrix. We define the $(n \times n)$−matrix-valued function $e^{At}$ as the following power series

$$e^{At} = I_n + At + A^2\frac{t^2}{2!} + A^3\frac{t^3}{3!} + A^4\frac{t^4}{4!} + \dots.$$

As with $e^{at}$, the radius of convergence is infinite. That is, the power series $e^{At}$ converges to a matrix for all $t < \infty$.

**Question 2.7:** Is $e^{At}$ actually a fundamental matrix for $(\frac{d}{dt} - A)\vec{x} = \vec{0}$?

**Answer:** Let's check! Recall that we need to check two things: are the columns of $e^{At}$ solutions to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ and are the columns linearly independent.

To check that the columns are solutions to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$, it suffices to show that $e^{At}$ is a solution to the corresponding matrix-valued ODE

$$
\begin{aligned}
\frac{d}{dt}e^{At} &= \frac{d}{dt}(I_n + At + A^2\frac{t^2}{2!} + A^3\frac{t^3}{3!} + A^4\frac{t^4}{4!} + ....) \\
&= \frac{d}{dt}I_n + \frac{d}{dt}At + \frac{d}{dt}A^2\frac{t^2}{2!} + \frac{d}{dt}A^3\frac{t^3}{3!} + \frac{d}{dt}A^4\frac{t^4}{4!} + .... \\
&= 0_n + A + A^2\frac{t}{1!} + A^3\frac{t^2}{2!} + A^4\frac{t^3}{3!} + .... \\
&= A((I_n + At + A^2\frac{t^2}{2!} + A^3\frac{t^3}{3!} + A^4\frac{t^4}{4!} + ....) \\
&= Ae^{At}
\end{aligned}
$$

Justifying the second line in the calculation above takes a bit of finesse, since we are switching the limit involved in the derivative with the limit of the sums. This requires that we check just *how* $e^{At}$ converges, but we will sweep these details under the rug. Suffice it to say that it works and thus, $e^{At}$ **is a solution to the corresponding matrix-valued ODE** $Z' = AZ$. This means that the columns of $e^{At}$ are solutions to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$.

Next, we need to check that the columns of $e^{At}$ are linearly independent. Our test for independence is to take the Wronskian of the columns. Recall that this is just the determinant of the matrix. By Abel's Theorem, we need only check at a single time, $t$. Let that time be $t = 0$. Hence,

$$
\begin{aligned}
W(0) &= det(e^{A \cdot 0}) \\
&= det(I_n + A0 + A^2\frac{0^2}{2!} + A^3\frac{0^3}{3!} + A^4\frac{0^4}{4!} + ....) \\
&= det(I_n) \\
&= 1
\end{aligned}
$$

Thus, the columns of $e^{At}$ are linearly independent solutions to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$. $e^{At}$ **is a fundamental matrix for** $(\frac{d}{dt} - A)\vec{x} = \vec{0}$**.**

It turns out, by the Existence and Uniqueness Theorem, that $e^{At}$ is the unique fundamental matrix which satisfies the matrix-valued Initial Value Problem

$$
Z' = AZ \quad , \quad Z(0) = I_n
$$

This is easily seen by applying our vector Existence and Uniqueness Theorem to the columns of $Z$.

## 2.3 How to find $e^{At}$ explicitly

Finding $e^{At}$ from the power series is almost always impossible. But, **Jordan Form makes it easy.**

**Change of Basis** Let $A$ be an $(n \times n)-$matrix. Let $J$ be the Jordan Form of $A$ and $S$ be the change of basis matrix so that $AS = SJ$. Given a vector-valued ODE,

$$\frac{d}{dt}\vec{x} = A\vec{x},$$

we want to change coordinates, i.e., change our basis, to the best coordinates, the ones in which the linear transformation represented by $A$ is simplest.

Doing the change of basis substitution, $\vec{x} = S\vec{y}$, we get the new system of ODEs $S\frac{d}{dt}\vec{y} = AS\vec{y}$. Rearranging, and noticing that by definition, $J = S^{-1}AS$, we have

$$\frac{d}{dt}\vec{y}(t) = J\vec{y}(t)$$

Now, we know that $e^{Jt}$ is a fundamental matrix for $\frac{d}{dt}\vec{y} = J\vec{y}$. So, changing our coordinates back, I claim that $Se^{Jt}$ **is a fundamental matrix for** $\frac{d}{dt}\vec{x}(t) = A\vec{x}(t)$**.**

**Proof:** What does it mean to be a fundamental matrix for $\frac{d}{dt}\vec{x}(t) = A\vec{x}(t)$? It means that the columns of the matrix-valued function $Se^{Jt}$ form a linearly independent collection of solutions to $\frac{d}{dt}\vec{x}(t) = J\vec{x}(t)$. To see that the columns are solutions, we check that $Se^{Jt}$ satisfies the matrix equation $Z' = AZ$.

$$
\begin{aligned}
\frac{d}{dt}(Se^{Jt}) &= S\frac{d}{dt}e^{Jt} \\
&= SJe^{Jt} \\
&= SJ(S^{-1}S)e^{Jt} \\
&= ASe^{Jt}
\end{aligned}
$$

To see that the columns are linearly independent, we check the Wronskian. As always, by Abel's Theorem, we only need check at one time, let that be $t = 0$.

$$
\begin{aligned}
W(0) &= det(Se^{J0}) \\
&= det(SI_n) \\
&= det(S) \\
&\neq 0
\end{aligned}
$$

where we know that $det(S) \neq 0$ because the columns of $S$ form a basis.

In practice, as has already been discussed, it is usually sufficient to find a fundamental matrix, ANY fundamental matrix. But, since this section is about finding $e^{At}$, let's actually find $e^{At}$. Recall that $e^{At}$ **is the unique matrix-valued function which satisfies the matrix Initial Value Problem**

$$Z'(t) = AZ(t) \quad , \quad Z(0) = I_n$$

If we consider $Se^{Jt}S^{-1}$, then, clearly $Se^{J0}S^{-1} = I_n$. And,

$$
\begin{aligned}
\frac{d}{dt}(Se^{Jt}S^{-1}) &= S\frac{d}{dt}e^{Jt}S^{-1} \\
&= SJe^{Jt}S^{-1} \\
&= SJ(S^{-1}S)e^{Jt}S^{-1} \\
&= ASe^{Jt}S^{-1}.
\end{aligned}
$$

**Thus, $Se^{Jt}S^{-1} = e^{At}$. This means we only need to find $e^{Jt}$. Because $J$ has such a simple form, this is actually relatively easy.**

**How to find $e^{Jt}$:** We derive $e^{Jt}$ case-by-case and from the definition. Recall that by definition

$$e^{Jt} = I_n + Jt + J^2 \frac{t^2}{2!} + J^3 \frac{t^3}{3!} + J^4 \frac{t^4}{4!} + \dots.$$

1. If $J$ is of the form $J = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$, then

$$\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}^2 = \begin{pmatrix} \lambda_1^2 & 0 \\ 0 & \lambda_2^2 \end{pmatrix}$$

$$\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}^3 = \begin{pmatrix} \lambda_1^3 & 0 \\ 0 & \lambda_2^3 \end{pmatrix}$$

$$\vdots$$

$$\begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}^n = \begin{pmatrix} \lambda_1^n & 0 \\ 0 & \lambda_2^n \end{pmatrix}$$

Now, we just sum up the series. Adding component by component, we get

$$
\begin{aligned}
e^{Jt} &= I_n + Jt + J^2 \frac{t^2}{2!} + J^3 \frac{t^3}{3!} + J^4 \frac{t^4}{4!} + \dots. \\
&= \begin{pmatrix} 1 + \lambda_1 t + \lambda_1^2 \frac{t^2}{2!} + \dots & 0 \\ 0 & 1 + \lambda_2 t + \lambda_2^2 \frac{t^2}{2!} + \dots \end{pmatrix} \\
&= \begin{pmatrix} e^{\lambda_1 t} & 0 \\ 0 & e^{\lambda_2 t} \end{pmatrix}.
\end{aligned}
$$

2. If $J$ is of the form $J = \begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix}$, then

$$\begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix}^2 = \begin{pmatrix} \lambda^2 & 2\lambda \\ 0 & \lambda^2 \end{pmatrix}$$

$$\begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix}^3 = \begin{pmatrix} \lambda^3 & 3\lambda^2 \\ 0 & \lambda^3 \end{pmatrix}$$

$$\vdots$$

$$\begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix}^n = \begin{pmatrix} \lambda^n & n\lambda^{n-1} \\ 0 & \lambda^n \end{pmatrix}$$

Now, we just sum up the series. Adding component by component, we get

$$
\begin{aligned}
e^{Jt} &= I_n + Jt + J^2\tfrac{t^2}{2!} + J^3\tfrac{t^3}{3!} + J^4\tfrac{t^4}{4!} + \dots \\
&= \begin{pmatrix} 1 + \lambda t + \lambda^2\tfrac{t^2}{2!} + \dots & 0 + 1t + 2\lambda\tfrac{t^2}{2!} + 3\lambda^2\tfrac{t^3}{3!} + \dots \\ 0 & 1 + \lambda t + \lambda^2\tfrac{t^2}{2!} + \dots \end{pmatrix} \\
&= \begin{pmatrix} e^{\lambda t} & te^{\lambda t} \\ 0 & e^{\lambda t} \end{pmatrix}.
\end{aligned}
$$

3. If $J$ is of the form $J = \begin{pmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{pmatrix}$, then

$$
\begin{pmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{pmatrix}^2 = \begin{pmatrix} \lambda^2 & 2\lambda & 1 \\ 0 & \lambda^2 & 2\lambda \\ 0 & 0 & \lambda^2 \end{pmatrix}
$$

$$
\begin{pmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{pmatrix}^3 = \begin{pmatrix} \lambda^3 & 3\lambda^2 & 3\lambda \\ 0 & \lambda^3 & 3\lambda^2 \\ 0 & 0 & \lambda^3 \end{pmatrix}
$$

$$
\vdots
$$

$$
\begin{pmatrix} \lambda & 1 & 0 \\ 0 & \lambda & 1 \\ 0 & 0 & \lambda \end{pmatrix}^n = \begin{pmatrix} \lambda^n & n\lambda^{n-1} & \tfrac{n(n-1)}{2}\lambda^{n-2} \\ 0 & \lambda^n & n\lambda^{n-1} \\ 0 & 0 & \lambda^n \end{pmatrix}
$$

Again, we sum up the series, adding component by component.

$$
\begin{aligned}
e^{Jt} &= \begin{pmatrix} 1 + \lambda t + \lambda^2\tfrac{t^2}{2!} + \dots & 0 + 1t + 2\lambda\tfrac{t^2}{2!} + 3\lambda^2\tfrac{t^3}{3!} + \dots & 0 + 0 + 1\tfrac{t^2}{2!} + 3\lambda\tfrac{t^3}{3!} + \dots \\ 0 & 1 + \lambda t + \lambda^2\tfrac{t^2}{2!} + \dots & 0 + 1t + 2\lambda\tfrac{t^2}{2!} + 3\lambda^2\tfrac{t^3}{3!} + \dots \\ 0 & 0 & 1 + \lambda t + \lambda^2\tfrac{t^2}{2!} + \dots \end{pmatrix} \\
&= \begin{pmatrix} e^{\lambda t} & te^{\lambda t} & \tfrac{t^2}{2!}e^{\lambda t} \\ 0 & e^{\lambda t} & te^{\lambda t} \\ 0 & 0 & e^{\lambda t} \end{pmatrix}.
\end{aligned}
$$

4. The other cases are easily dealt with once we recall a fact about block diagonal matrices. If $A$ and $B$ are submatrices of a block diagonal matrix, then,

$$
\begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}^2 = \begin{pmatrix} A^2 & 0 \\ 0 & B^2 \end{pmatrix}
$$

Applying this rule to the other possible Jordan Forms easily gives their formulae. For instance, it shows that if a matrix, $M$, is block diagonal (all Jordan forms are block diagonal), then if $M = \begin{pmatrix} A & 0 \\ 0 & B \end{pmatrix}$ for some submatrices, $A, B$,

$$
e^{Mt} = \begin{pmatrix} e^{At} & 0 \\ 0 & e^{Bt} \end{pmatrix}
$$

However, **this only works with Block diagonal matrices.**

**Thus, to find $e^{At}$,**

- **Step 1:** Find $J$, the Jordan Form of $A$, and $S$ the change of basis matrix such that $AS = SJ$.

- **Step 2:** Use the formulas to write down $e^{Jt}$.

- **Step 3:** $e^{At} = Se^{Jt}S^{-1}$.

**Example 2.2:** Let $A = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}$. Give the matrix-valued function $e^{At}$.

By previous example, we know that $J = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ and $S = \begin{pmatrix} 1 & 0 \\ 0 & 1/2 \end{pmatrix}$.

Thus, $e^{Jt} = \begin{pmatrix} e^t & te^t \\ 0 & e^t \end{pmatrix}$ and

$$e^{At} = \begin{pmatrix} 1 & 0 \\ 0 & 1/2 \end{pmatrix} \begin{pmatrix} e^t & te^t \\ 0 & e^t \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

**Reflection Questions**

1. What is the series definition of $e^{At}$?

2. What type of object is $e^{At}$? Is it a matrix, vector, scalar-valued function, vector-valued function, or matrix-valued function?

3. What properties does $e^{At}$ have?

4. Given a $(2 \times 2)$ or $(3 \times 3)$ matrix, $A$, by what steps do we find $e^{At}$?

5. If $J = \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 1 \\ 0 & 0 & \lambda_2 \end{pmatrix}$ what is $e^{Jt}$?

6. If $J = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_1 & 0 & 0 \\ 0 & 0 & \lambda_2 & 1 \\ 0 & 0 & 0 & \lambda_2 \end{pmatrix}$ what is $e^{Jt}$?

7. If $J = \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 1 \\ 0 & 0 & \lambda_2 \end{pmatrix}$ what is $\frac{d}{dt}e^{Jt}$?

## 2.4 Solving $\frac{d}{dt}\vec{x} = A\vec{x}$

The whole point of our approach to this problem has been that every solution to $\frac{d}{dt}\vec{x} = A\vec{x}$ can be written as a linear combination of the columns of a fundamental matrix. If $J$ is the Jordan Form of $A$ and $S$ is the change of basis matrix such that $AS = SJ$, then $Se^{Jt}$ is a fundamental matrix for the system $\frac{d}{dt}\vec{x} = A\vec{x}$. Thus, every solution of $\frac{d}{dt}\vec{x} = A\vec{x}$ can be written as

$$\vec{x}(t) = Se^{Jt}\vec{c}$$

for some vector $\vec{c} \in \mathbb{R}^n$.

**Definition 2.9.** The **general solution** to a system of equations is an expression involving parameters, $c_1, c_2, ..., c_n$, such that two things hold:

1. For every choice of parameters, $c_1, c_2, ..., c_n$, the expression is a solution to the given system of equations.

2. Every solution to the given system of equations can be written as some particular choice of parameters in the expression.

Not every system of equations has a general solution. In fact, general $1^{st}$ order linear system of equations, $\frac{d}{dt}\vec{x}(t) = A(t)\vec{x}(t)$ may not have a general solution. However, if we restrict to the constant coefficient case, then $\frac{d}{dt}\vec{x} = A\vec{x}$ does have a general solution, and the general solution is

$$\vec{x}(t) = Se^{Jt}\vec{c}.$$

**Example 2.3.** Let $A = \begin{pmatrix} 3 & -18 \\ 2 & -9 \end{pmatrix}$. Give the general solution to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$.

Since the general solution is $Se^{Jt}\vec{c}$, we need to find $S$, $J$ and then $e^{Jt}$.

1. First, we find $J$ and $S$. Checking the characteristic equation, we find that

$$det(A - \lambda I) = (\lambda + 3)^2 = 0.$$

Hence, we have an eigenvalue, $\lambda = -3$, with algebraic multiplicity $m = 2$. To find eigenvectors, we solve

$$(A - (-3)\lambda)\vec{x} = \begin{pmatrix} 6 & -18 \\ 2 & -6 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \vec{0}$$

to get a single linearly independent eigenvector, $\xi = \begin{pmatrix} 3 \\ 1 \end{pmatrix}$. To build our Jordan chain, we then solve

$$(A - (-3)\lambda)\vec{x} = \begin{pmatrix} 6 & -18 \\ 2 & -6 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \end{pmatrix}$$

to get the generalized eigenvector $\vec{\eta} = \begin{pmatrix} 1/2 \\ 0 \end{pmatrix}$.

Thus, $J = \begin{pmatrix} -3 & 1 \\ 0 & -3 \end{pmatrix}$ and $S = \begin{pmatrix} 3 & 1/2 \\ 1 & 0 \end{pmatrix}$.

2. Now, we write down $e^{Jt}$

$$e^{Jt} = \begin{pmatrix} e^{-3t} & te^{-3t} \\ 0 & e^{-3t} \end{pmatrix}$$

3. Finally, we can write down the general solution.

$$Se^{Jt}\vec{c} = \begin{pmatrix} 3 & 1/2 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} e^{-3t} & te^{-3t} \\ 0 & e^{-3t} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}$$

Sometimes it is helpful to actually do the matrix multiplication. In that case the general solution becomes

$$Se^{Jt}\vec{c} = c_1 \begin{pmatrix} 3 \\ 1 \end{pmatrix} e^{-3t} + c_2(\begin{pmatrix} 3 \\ 1 \end{pmatrix} te^{-3t} + \begin{pmatrix} 1/2 \\ 0 \end{pmatrix} e^{-3t})$$

**Example 2.4.** Let $A = \begin{pmatrix} 1 & 2 \\ -3 & 0 \end{pmatrix}$. Give the general solution to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$.

1. FIrst, find $J$ and $S$. Checking the Characteristic equation, we find that

$$det(A - \lambda I) = \lambda^2 - \lambda + 6 = 0.$$

Using the quadratic equation, we see that $\lambda = 1/2 \pm i\sqrt{23}/2$. To find eigenvectors, we first solve

$$(A - (1/2 + i\sqrt{23}/2)\lambda)\vec{x} = \begin{pmatrix} 1/2 - i\sqrt{23}/2 & 2 \\ -3 & -1/2 - i\sqrt{23}/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \vec{0}$$

to get the eigenvector, $\xi_1 = \begin{pmatrix} -1/6 - i\sqrt{23}/6 \\ 1 \end{pmatrix}$.

Next, we pick the other eigenvalue and solve

$$(A - (1/2 - i\sqrt{23}/2)\lambda)\vec{x} = \begin{pmatrix} 1/2 + i\sqrt{23}/2 & 2 \\ -3 & -1/2 + i\sqrt{23}/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \vec{0}$$

to get the eigenvector, $\xi_2 = \begin{pmatrix} -1/6 + i\sqrt{23}/6 \\ 1 \end{pmatrix}$.

Thus, $J = \begin{pmatrix} 1/2 + i\sqrt{23}/2 & 0 \\ 0 & 1/2 - i\sqrt{23}/2 \end{pmatrix}$ and $S = \begin{pmatrix} -1/6 - i\sqrt{23}/6 & -1/6 + i\sqrt{23}/6 \\ 1 & 1 \end{pmatrix}$.

39

2. Now, we need to write down $e^{Jt}$

$$e^{Jt} = \begin{pmatrix} e^{(1/2+i\sqrt{23}/2)t} & 0 \\ 0 & e^{(1/2-i\sqrt{23}/2)t} \end{pmatrix}$$

3. Finally, we can write down the general solution.

$$Se^{Jt}\vec{c} = \begin{pmatrix} -1/6 - i\sqrt{23}/6 & -1/6 + i\sqrt{23}/6 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} e^{(1/2+i\sqrt{23}/2)t} & 0 \\ 0 & e^{(1/2-i\sqrt{23}/2)t} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}$$

Again, it can be helpful to do the matrix multiplication. In that case the general solution becomes

$$Se^{Jt}\vec{c} = c_1 \begin{pmatrix} -1/6 - i\sqrt{23}/6 \\ 1 \end{pmatrix} e^{(1/2+i\sqrt{23}/2)t} + c_2 \begin{pmatrix} -1/6 + i\sqrt{23}/6 \\ 1 \end{pmatrix} e^{(1/2-i\sqrt{23}/2)t}.$$

**Question 2.9:** What if we want real-valued solutions?

**Answer:** There is nothing wrong with complex-valued solutions! But, if you insist on having real-valued solutions, that isn't too hard. Let's review a few facts about complex numbers.

First, recall that a complex number, $z \in \mathbb{C}$, can be written as $z = a + ib$ for two real numbers $a, b$. The number $a$ is called the **real part** of $z$, and the number $b$ is called the **imaginary part** of $z$. If $\vec{x}(t)$ is a complex, vector-valued solution, for each time, $t$, $\vec{x}(t)$ is a vector with complex entries. Since we can split each complex number into real and imaginary parts, we can split complex vectors into real and imaginary parts. This lets us we can split $\vec{x}(t)$ into real and imaginary parts. That is, we can write

$$\vec{x}(t) = \vec{u}(t) + i\vec{v}(t)$$

where $\vec{u}(t), \vec{v}(t)$ are real, vector-valued functions.

Next we observe three simple facts about complex numbers: If $a, b \in \mathbb{R}$

- The product of two real numbers is real. $ab \in \mathbb{R}$.

- The product of a real number and an imaginary number is imaginary. $a(ib) = i(ab)$.

- Two complex numbers are equal if and only if their real parts and their imaginary parts are equal. $a + ib = c + id$ iff $a + c$ and $b = d$.

Thus, because $A$ only has real values, if $\vec{x}(t) = \vec{u}(t) + i\vec{v}(t)$ is a solution to $\frac{d}{dt}\vec{x} = A\vec{x}$, then

$$\frac{d}{dt}(\vec{u}(t) + i\vec{v}(t)) = A(\vec{u}(t) + i\vec{v}(t))$$

and therefore,

$$\frac{d}{dt}\vec{u}(t) + i\frac{d}{dt}\vec{v}(t) = A\vec{u}(t) + iA\vec{v}(t).$$

Thus, because real parts and imaginary parts must be equal, we have that

$$\begin{array}{rcl} \frac{d}{dt}\vec{u}(t) & = & A\vec{u}(t) \\ \frac{d}{dt}\vec{v}(t) & = & A\vec{v}(t) \end{array}$$

So, **if we have a complex-valued solution, $\vec{x}(t) = \vec{u}(t) + i\vec{v}(t)$, and $A$ is real-valued, then the real and complex parts, $\vec{u}(t)$ and $\vec{v}(t)$, are real-valued solutions.**

In the example above, that means that we need to take one of the solutions and slit it into real and imaginary parts. We do that using Euler's Formula.

$$
\begin{aligned}
\begin{pmatrix} -1/6 + i\sqrt{23}/6 \\ 1 \end{pmatrix} e^{(1/2 - i\sqrt{23}/2)t} &= \left( \begin{pmatrix} -1/6 \\ 1 \end{pmatrix} + i \begin{pmatrix} \sqrt{23}/6 \\ 0 \end{pmatrix} \right) e^{(1/2 - i\sqrt{23}/2)t} \\
&= \left( \begin{pmatrix} -1/6 \\ 1 \end{pmatrix} + i \begin{pmatrix} \sqrt{23}/6 \\ 0 \end{pmatrix} \right) e^{\frac{1}{2}t} (\cos(\tfrac{-\sqrt{23}}{2}t) + i\sin(\tfrac{-\sqrt{23}}{2}t)) \\
&= \left( \begin{pmatrix} -1/6 \\ 1 \end{pmatrix} e^{\frac{1}{2}t} \cos(\tfrac{-\sqrt{23}}{2}t) - \begin{pmatrix} \sqrt{23}/6 \\ 0 \end{pmatrix} e^{\frac{1}{2}t} \sin(\tfrac{-\sqrt{23}}{2}t) \right) \\
&\quad + i\left( \begin{pmatrix} -1/6 \\ 1 \end{pmatrix} e^{\frac{1}{2}t} \sin(\tfrac{-\sqrt{23}}{2}t) + \begin{pmatrix} \sqrt{23}/6 \\ 0 \end{pmatrix} e^{\frac{1}{2}t} \cos(\tfrac{-\sqrt{23}}{2}t) \right) \\
&= \vec{u}(t) + i\vec{v}(t)
\end{aligned}
$$

**Question 2.10:** Are $\vec{u}(t)$ and $\vec{v}(t)$ linearly independent?

**Answer:** Let's check! As always, Abel's theorem allows us to choose a single time, $t$, to check at. Let's choose $t = 0$.

$$
\begin{aligned}
W[\vec{u}, \vec{v}](0) &= det\begin{pmatrix} \frac{1}{6} & \frac{\sqrt{23}}{6} \\ 1 & 0 \end{pmatrix} \\
&= det(S) \\
&\neq 0
\end{aligned}
$$

Thus, **they are always linearly independent.**

**Example 2.5.** Let $A = \begin{pmatrix} 5 & -4 & 0 \\ 1 & 0 & 2 \\ 0 & 2 & 5 \end{pmatrix}$. Give the general solution to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$.

1. Again, we find $J$ and $S$. Checking the characteristic equation, we find that

$$det(A - \lambda I) = (\lambda - 5)^2(-\lambda) = 0.$$

Thus, our eigenvalues are $\lambda = 5$ and $\lambda = 0$. To find eigenvectors, we first solve

$$(A - (5)\lambda)\vec{x} = \begin{pmatrix} 0 & -4 & 0 \\ 1 & -5 & 2 \\ 0 & 2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \vec{0}$$

to get the eigenvector, $\xi_1 = \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix}$.

Next, we pick the other eigenvalue and solve

$$(A - (0)\lambda)\vec{x} = \begin{pmatrix} 5 & -4 & 0 \\ 1 & 0 & 2 \\ 0 & 2 & 5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \vec{0}$$

to get the eigenvector, $\xi_2 = \begin{pmatrix} -2 \\ -5/2 \\ 1 \end{pmatrix}$.

Now, we build our Jordan chain off of $\xi_1 = \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix}$. We solve

$$(A - (5)\lambda)\vec{x} = \begin{pmatrix} 0 & -4 & 0 \\ 1 & -5 & 2 \\ 0 & 2 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \xi_1 = \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix}$$

to get the generalized eigenvector, $\vec{\eta} = \begin{pmatrix} 5/2 \\ 1/2 \\ 0 \end{pmatrix}$.

Thus, $J = \begin{pmatrix} 5 & 1 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 0 \end{pmatrix}$ and $S = \begin{pmatrix} -2 & 5/2 & -2 \\ 0 & 1/2 & -5/2 \\ 1 & 0 & 1 \end{pmatrix}$.

2. Now, we need to write down $e^{Jt}$

$$e^{Jt} = \begin{pmatrix} e^{5t} & te^{5t} & 0 \\ 0 & e^{5t} & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

3. Finally, we can write down the general solution.

$$Se^{Jt}\vec{c} = \begin{pmatrix} -2 & 5/2 & -2 \\ 0 & 1/2 & -5/2 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} e^{5t} & te^{5t} & 0 \\ 0 & e^{5t} & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix}$$

Again, it can be helpful to do the matrix multiplication. In that case the general solution becomes

$$\vec{x}(t) = c_1 \begin{pmatrix} -2 \\ 0 \\ 1 \end{pmatrix} e^{5t} + c_2 (\begin{pmatrix} -2 \\ 0 \\ 1 \end{pmatrix} te^{5t} + \begin{pmatrix} 5/2 \\ 1/2 \\ 0 \end{pmatrix} e^{5t}) + c_3 \begin{pmatrix} -2 \\ -5/2 \\ 1 \end{pmatrix}$$

## 2.5 How do solutions behave?

Now that we know how to find a basis for solutions to equations of the form $(\frac{d}{dt} - A)\vec{x} = \vec{0}$, we might ask, how do they behave? What happens to these solutions as time goes on? We restrict ourselves to considering the $2-$dimensional case.

**If $A$ is a $(2 \times 2)-$matrix, then there are only two possibilities.** Either $A$ is diagonalizable, in which case the general solution looks like

$$\vec{x}(t) = c_1\vec{\xi_1}e^{\lambda_1 t} + c_2\vec{\xi_2}e^{\lambda_2 t}$$

or, $A$ is defective, in which case the general solution looks like

$$\vec{x}(t) = c_1\vec{\xi}e^{\lambda t} + c_2(\vec{\xi}te^{\lambda t} + \vec{\eta}e^{\lambda t}).$$

In either case, every solution is a linear combination of functions that look like $\vec{\xi}e^{\lambda t}$ or $\vec{\xi}te^{\lambda t}$ for some choice of $\lambda$ and $\vec{\xi}$.

**Question 2.11:** What kind of objects ARE $\vec{\xi}e^{\lambda t}$ and $\vec{\xi}te^{\lambda t}$?

**Answer:** These are vector-valued functions. Recall that this means that for every time $t$, $\vec{\xi}e^{\lambda t}$ **is a vector in $\mathbb{R}^2$. As $t$ changes, $\vec{\xi}e^{\lambda t}$ gives different vectors in $\mathbb{R}^2$. Thus, the function $\vec{\xi}e^{\lambda t}$ describes a path in $\mathbb{R}^2$ parametrized by $t$.**

**Definition 2.10.** An **integral curve** of a system of ODEs, $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$, is a curve which traces the path of a solution. That is, for all the vectors, $\vec{x}, \vec{y}$, on the integral curve, there is a solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$ which goes through both $\vec{x}$ and $\vec{y}$.

This definition captures the notion that solutions describe paths in $\mathbb{R}^2$.

**Question 2.12:** How do $\vec{\xi}e^{\lambda t}$ and $\vec{\xi}te^{\lambda t}$ behave as time goes on?

**Answer:** Obviously this depends upon $\vec{\xi}$ and $\lambda$. But let's concentrate on what solutions do as $t \to \infty$. This is a simple limit calculation.

$$\begin{aligned} lim_{t\to\infty}|\vec{\xi}e^{\lambda t}| &= |\vec{\xi}|lim_{t\to\infty}|e^{\lambda t}| \\ &= \begin{cases} 0 & \text{if } Re(\lambda) < 0 \\ \infty & \text{if } Re(\lambda) > 0 \\ |\vec{\xi}| & \text{if } Re(\lambda) = 0 \end{cases} \end{aligned}$$

Similarly,

$$
\begin{aligned}
lim_{t \to \infty} |\vec{\xi} t e^{\lambda t}| &= |\vec{\xi}| lim_{t \to \infty} |t e^{\lambda t}| \\
&= \begin{cases} 0 & \text{if } Re(\lambda) < 0 \\ \infty & \text{if } Re(\lambda) > 0 \\ \infty & \text{if } Re(\lambda) = 0 \end{cases}
\end{aligned}
$$

As you can see, this means there are three types of solutions: solutions which eventually "blow up" by going off to infinity, solutions which asymptotically approach the origin, and solutions which neither blow up nor asymptotically approach the origin. **These completely determine the behavior of solutions.**

**Definition 2.11.** A solution to $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ is called an **equilibrium solution** or a **critical point** of $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ if $\frac{d}{dt}\vec{x}(t) = \vec{0}$. That is, an equilibrium solution is one which is constant.

Observe that for homogeneous $1^{st}$ order linear systems of ODEs, the zero vector, $\vec{x}(t) \equiv \vec{0}$, is always an equilibrium solution.

**Classifying Critical Points** To classify critical points, we care about how solutions near a critical point behave. From the discussion above, we can see that solutions only behave one of three ways, we classify along those lines.

**Definition 2.12.** A critical point of $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ is called **unstable** if ANY solution nearby diverges to infinity.

**Definition 2.13.** A critical point of $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ is called **asymptotically stable** if ALL solution nearby asymptotically approach that critical point.

**Definition 2.14.** A critical point of $(\frac{d}{dt} - A)\vec{x} = \vec{0}$ is called **stable** if ALL solutions nearby neither diverge to infinity nor asymptotically approach that critical point.

It is clear from the previous page that $\vec{0}$ is an **unstable critical point if any of the eigenvalues of $A$ have positive real part.** Similarly, $\vec{0}$ is an **asymptotically stable**

**critical point if al of the eigenvalues of $A$ have negative real part.** And, finally, $\vec{0}$ is a **stable critical point only if all of the eigenvalues of $A$ are non-zero, but have real part equal to zero.**

Note that this does not completely classify all critical points. If $\lambda = 0$ is an eigenvalue, then all eigenvectors, $\vec{\xi}$, associated to $\lambda = 0$, will be critical points of $(\frac{d}{dt} - A)\vec{x} = \vec{0}$. But, if $A$ has another non-zero eigenvalue with negative real part, then these points will not necessarily be asymptotically stable, stable, or unstable.

**Question:** How do we draw integral curves?

**Answer:** We draw integral curves to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$ by first writing down the general solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$. There are only 2 possibilities for a $(2 \times 2)-$matrix.
Either $A$ is diagonalizable, in which case the general solution looks like

$$\vec{x}(t) = c_1\vec{\xi_1}e^{\lambda_1 t} + c_2\vec{\xi_2}e^{\lambda_2 t}$$

or, $A$ is defective, in which case the general solution looks like

$$\vec{x}(t) = c_1\vec{\xi}e^{\lambda t} + c_2(\vec{\xi}te^{\lambda t} + \vec{\eta}e^{\lambda t}).$$

By the Existence and Uniqueness theorem, for every vector $\vec{x}_0 \in \mathbb{R}^2$, there exists a solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$ which satisfies the initial condition, $\vec{x}(0) = \vec{x}_0$. Plugging in the initial conditions, we see that the constants $c_1, c_2$ simply determine the linear combinations of $\vec{\xi_1}, \vec{\xi_2}$ or $\vec{\xi}, \vec{\eta}$ which sum up to $\vec{x}_0$.
From this point, since we now have an explicit formula for $\vec{x}(t)$, we need only take the limit,

$$lim_{t\to\infty}\vec{x}(t)$$

and observe the relative speeds with which $e^{\lambda t}$ and $te^{\lambda t}$ grow or decay.

**Reflection Questions**

1. How might we find critical points of $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$?

2. What is an integral curve? Is it possible for different solutions to lie on the same integral curve?

3. What is the general process for drawing integral curves of $(\frac{d}{dt} - A)\vec{x}(t) = \vec{0}$?

4. Suppose $A$ has two eigenvalues, $\lambda = 0$ and $\lambda = 24$. How many critical points will $A$ have? Will they be asymptotically stable, stable, or unstable?

5. Suppose $A$ is a $(2 \times 2)-$matrix which has one eigenvalue, $\lambda = -3$. How many critical points will $A$ have? Will they be asymptotically stable, stable, or unstable?

6. If $\lambda_1 > \lambda_2$, which function grows faster as $t \to \infty$, $e^{\lambda_1 t}$ or $e^{\lambda_2 t}$?

7. If $\lambda_1 > \lambda_2$, which function shrinks faster as $t \to -\infty$, $e^{\lambda_1 t}$ or $e^{\lambda_2 t}$?

8. Which function grows faster as $t \to \infty$, $e^{\lambda t}$ or $te^{\lambda t}$?

9. Which function shrinks faster as $t \to -\infty$, $e^{\lambda t}$ or $te^{\lambda t}$?

## 2.6 Non-homogeneous Equations

Now that we know a lot about $1^{st}$ order linear systems of homogeneous ODEs, we consider a closely related class of systems, $1^{st}$ order linear systems of non-homogeneous ODEs. In general, these systems are written as

$$\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + \vec{g}(t).$$

Where $\vec{x}(t)$ and $\vec{g}(t)$ are **vector-valued functions**, i.e., vectors whose entries are functions, and $A$ is a square matrix with constant entries.

**Question 2.13:** What does the set of solutions to $\frac{d}{dt}\vec{x}(t) = A\vec{x}(t) + \vec{g}(t)$ look like? What kind of structure does it have?

**Answer:** Just as before, we re-write the equation $\frac{d}{dt}\vec{x} = A\vec{x} + \vec{g}$ as the following:

$$(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t).$$

Except that now, because our equations are non-homogeneous, solutions are no longer in the kernel of the operator $(\frac{d}{dt} - A)$. This changes everything.

**Solutions to non-homogeneous equations do not form a vector space.** To see this, let $\vec{x}(t)$ and $\vec{y}(t)$ be solutions to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$.

$$\begin{aligned}
(\frac{d}{dt} - A)(\vec{x}(t) + \vec{y}(t)) &= (\frac{d}{dt} - A)\vec{x}(t) + (\frac{d}{dt} - A)\vec{y}(t) \\
&= \vec{g}(t) + \vec{g}(t) \\
&\neq \vec{g}(t).
\end{aligned}$$

However, **the difference between any two solutions to a non-homogeneous linear system of ODEs is a solution to the corresponding homogeneous equation.** That is, if we again let $\vec{x}(t)$ and $\vec{y}(t)$ be solutions to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$.

$$\begin{aligned}
(\frac{d}{dt} - A)(\vec{x}(t) - \vec{y}(t)) &= (\frac{d}{dt} - A)\vec{x}(t) - (\frac{d}{dt} - A)\vec{y}(t) \\
&= \vec{g}(t) - \vec{g}(t) \\
&= \vec{0}.
\end{aligned}$$

Since we know that every solution to the corresponding homogeneous linear systems of ODEs can be written as $Se^{Jt}\vec{c}$, we can write every solution to the non-homogeneous equation as $\vec{x}(t) + Se^{Jt}\vec{c}$. **Therefore, the general solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ is**

$$\vec{x}_p(t) + Se^{Jt}\vec{c}$$

where $\vec{x}_p$ is any solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$, and $J$ is the Jordan form of $A$ and $S$ is the change of basis matrix such that $AS = SJ$. Thus, the set of solutions to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ forms an **affine linear subspace.**

**Definition 2.15.** A set is called an **affine linear subspace** if it is the translate of a linear subspace. This means it takes the form $\vec{v} + V$ where $V$ is a linear subspace and $\vec{v} \neq \vec{0}$.

**Example:** For any $b \neq 0$, the set $\{(x, y) \in \mathbb{R}^2 | y = mx + b\}$ is an affine linear subspace of $\mathbb{R}^2$.

**This means that to find the general solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$, we need only find a single solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ and then add on the general solution to the corresponding homogeneous equation.**

**Question 2.14:** How do we find a solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$, then?

**Answer:** There are many ways to find a solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$. We will focus on two methods: Variation of Parameters and Change of Basis.

**Variation of Parameters**

Recall the method of Integrating Factors. For a $1^{st}$ Order linear non-homogeneous ODE,

$$x'(t) = a(t)x(t) + b(t)$$

we rearranged the equation like so

$$x'(t) - a(t)x(t) = b(t)$$

and then searched for a function, $\mu(t)$, such that $\mu(t)x'(t) - a(t)\mu(t)x(t) = (\mu(t)x(t))'$. Expanding this equation, we got

$$\mu(t)x'(t) - a(t)\mu(t)x(t) = \mu(t)x'(t) + \mu'(t)x(t)$$

which meant that $\mu'(t) = -a(t)\mu(t)$. Integrating, we saw that $\mu(t) = e^{-\int a(t)dt}$.

Thus, we could solve $x'(t) = a(t)x(t) + b(t)$ as follows,

$$
\begin{aligned}
x'(t) - a(t)x(t) &= b(t) \\
(\mu(t)x(t))' &= \mu(t)b(t) \\
\mu(t)x(t) &= \int \mu(t)b(t)dt + C \\
x(t) &= \mu(t)^{-1} \int \mu(t)b(t)dt + C\mu(t)^{-1}
\end{aligned}
$$

**Here is another way to look at the method of Integrating Factors.** We know that the general solution to the homogeneous equation $x'(t) = a(t)x(t)$ is given by

$$x(t) = e^{\int_0^t a(s)ds} c.$$

So what if the solution to the corresponding non-homogeneous equation $x'(t) = a(t)x(t) + b(t)$ is given by

$$x(t) = e^{\int_0^t a(s)ds} c(t).$$

for some function $c(t)$? Plugging in, we solve for $c(t)$. That is, assume $e^{\int_0^t a(s)ds} c(t)$ is a solution to $x'(t) = a(t)x(t) + b(t)$. Then

$$
\begin{aligned}
(e^{\int_0^t a(s)ds} c(t))' &= a(t)e^{\int_0^t a(s)ds} c(t) + b(t) \\
a(t)e^{\int_0^t a(s)ds} c(t) + e^{\int_0^t a(s)ds} c'(t) &= a(t)e^{\int_0^t a(s)ds} c(t) + b(t) \\
e^{\int_0^t a(s)ds} c'(t) &= b(t) \\
c'(t) &= e^{-\int_0^t a(s)ds} b(t) \\
c(t) &= \int e^{-\int a(r)dr} b(s)ds + C
\end{aligned}
$$

Thus, our solution to $x'(t) = a(t)x(t) + b(t)$ is given by

$$e^{\int_0^t a(s)ds} c(t) = e^{\int_0^t a(s)ds} \left( \int e^{-\int a(r)dr} b(s)ds + C \right).$$

Note that this is the same as we had derived previously.

If you recall, you did the same thing in 307 with $2^{nd}$ order linear ODEs. That is, if we have ah homogeneous equation $y'' + b(t)y' + c(t)y = 0$ with general solution $x(t) = c_1 y_1(t) + c_2 y_2(t)$, then we found a solution to the corresponding non-homogeneous equation, $y'' + b(t)y' + c(t)y = d(t)$ by assuming that the solution was of the form

$$x(t) = y_1(t)c_1(t) + y_2(t)c_2(t).$$

Again, actually solving for $c_1(t), c_2(t)$ is done by plugging into the equation. But we will not calculate them here. Instead, we give a more general solution.

If we have the homogeneous linear system of ODEs $(\frac{d}{dt} - A)\vec{x} = \vec{0}$, we know that for any fundamental matrix $\Phi(t)$, the general solution can be written as

$$\vec{x}(t) = \Phi(t)\vec{c}$$

So, maybe a solution to the corresponding non-homogeneous system $(\frac{d}{dt} - A)\vec{x} = \vec{g}(t)$ is of the form

$$\vec{x}(t) = \Phi(t)\vec{c}(t).$$

We find the vector-valued function $\vec{c}(t)$ by assuming that $\Phi(t)\vec{c}(t)$ is a solution and plugging it it.

$$\begin{aligned}
(\tfrac{d}{dt} - A)\Phi(t)\vec{c}(t) &= \vec{g}(t) \\
A\Phi(t)\vec{c}(t) + \Phi(t)\tfrac{d}{dt}\vec{c}(t) - A\Phi(t)\vec{c}(t) &= \vec{g}(t) \\
\Phi(t)\tfrac{d}{dt}\vec{c}(t) &= \vec{g}(t) \\
\tfrac{d}{dt}\vec{c}(t) &= \Phi(t)^{-1}\vec{g}(t)
\end{aligned}$$

thus, integrating, we see that

$$\vec{c}(t) = \int_0^t \Phi^{-1}(s)g(s)ds + \vec{c}$$

**Hence, we have a solution to** $(\tfrac{d}{dt} - A)\vec{x} = \vec{g}(t)$ **given by**

$$\vec{x}(t) = \Phi(t)\vec{c}(t) = \Phi(t)\int_0^t \Phi^{-1}(s)g(s)ds + \Phi(t)\vec{c}.$$

On closer inspection, **this is clearly the general solution, since it is written as a particular solution plus the general solution to the corresponding homogeneous equation.**

Observe that while this holds for ANY fundamental matrix, $\Phi(t)$, we can write is very suggestively for the fundamental matrix, $e^{At}$.

$$\vec{x}(t) = e^{At}\int_0^t e^{-As}g(s)ds + e^{At}\vec{c}.$$

This method of getting the general solution to $(\tfrac{d}{dt} - A)\vec{x} = \vec{g}(t)$ is very convenient for analysis, but rather cumbersome for calculation. For calculating explicit solutions, it is often easier to use a Change of Basis.

### Change of Basis

Just as we have throughout this course, we wish to change coordinates to find the simplest representation of $(\tfrac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$. As we have seen, this means finding the Jordan Form of $A$.

So, if $J$ is the Jordan form of $A$ and $S$ the change of basis matrix such that $AS = SJ$, we change variables by letting $\vec{x} = S\vec{y}$. Thus, we have the new equation:

$$\begin{aligned}
\tfrac{d}{dt}S\vec{y}(t) &= AS\vec{y}(t) + \vec{g}(t) \\
\tfrac{d}{dt}\vec{y}(t) &= J\vec{y}(t) + S^{-1}\vec{g}(t) \\
\tfrac{d}{dt}\vec{y}(t) &= J\vec{y}(t) + \vec{h}(t)
\end{aligned}$$

Case 1. If $A$ is diagonalizable, then $J$ is diagonal, and this reduces to the system of equations

$$\begin{aligned}
y_1'(t) &= \lambda_1 y_1(t) + h_1(t) \\
y_2'(t) &= \lambda_2 y_2(t) + h_2(t) \\
y_3'(t) &= \lambda_3 y_3(t) + h_3(t) \\
&\vdots \\
y_n'(t) &= \lambda_n y_n(t) + h_n(t)
\end{aligned}$$

We recall from 307, that we can solve $1^{st}$ order linear ODEs using the method of Integrating Factors (which is reviewed above in the Variation of Parameters section). Thus, the solution to $\frac{d}{dt}\vec{y}(t) = J\vec{y}(t) + \vec{h}(t)$ is a vector-valued function, $\vec{y}(t)$ which has entries

$$y_i(t) = e^{\lambda_i t} \int_0^t e^{-\lambda_i s} h_i(s) ds + c_i e^{\lambda_i t}.$$

Our solution to $\frac{d}{dt}\vec{x}(t) = A\vec{y}(t) + \vec{g}(t)$, then, is given by changing coordinates back by $\vec{x} = S\vec{y}$.

Case 2. If $J$ is not diagonal, then the system of equations may look something like this

$$
\begin{array}{rlll}
y_1'(t) & = & \lambda_1 y_1(t) \;+ y_2(t) & +h_1(t) \\
y_2'(t) & = & \lambda_1 y_2(t) & +h_2(t) \\
y_3'(t) & = & \lambda_3 y_3(t) & +h_3(t) \\
& \vdots & & \\
y_n'(t) & = & \lambda_n y_n(t) \;+ h_n(t) &
\end{array}
$$

Notice that while some equations may involve more that one component function $y_i(t)$, **the last row will always only involve** $y_n(t)$. Thus, we can solve the last row using Integrating Factors, get a solution and substitute that concrete solution up into the next equation. In this way we can successively solve all of the equations from the bottom to the top using Integrating Factors.

**Example 2.6 .** Let $A = \begin{pmatrix} 5 & -4 & 0 \\ 1 & 0 & 2 \\ 0 & 2 & 5 \end{pmatrix}$. Let $\vec{g}(t) = \begin{pmatrix} 0 \\ 0 \\ e^{-t} \end{pmatrix}$. Give the general solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$.

We first change bases. By previous calculation, we know that $J = \begin{pmatrix} 5 & 1 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 0 \end{pmatrix}$ and $S = \begin{pmatrix} 2 & 5/2 & -2 \\ 0 & 1/2 & -5/2 \\ 1 & 0 & 1 \end{pmatrix}$.

Calculating $S^{-1}\vec{g}(t)$, we get that $\vec{h}(t) = \begin{pmatrix} \frac{21}{17}e^{-t} \\ \frac{-20}{17}e^{-t} \\ \frac{-4}{17}e^{-t} \end{pmatrix}$.

Thus, we need to solve the system of equations:

$$
\begin{array}{rlll}
y_1'(t) & = & 5y_1(t) \;+ y_2(t) & +\frac{21}{17}e^{-t} \\
y_2'(t) & = & 5y_2(t) & +\frac{-20}{17}e^{-t} \\
y_3'(t) & = & 0y_3(t) \;+ \frac{-4}{17}e^{-t} &
\end{array}
$$

We solve from the bottom to the top. In this case, we can simply integrate the equation, $y_3'(t) = \frac{-4}{17}e^{-t}$, to get that

$$y_3(t) = \frac{4}{17}e^{-t} + c_3$$

To solve for $y_2(t)$ we use Integrating Factors. By multiplying by the integrating factor, $e^{-5t}$, have that

$$
\begin{aligned}
y_2'(t) &= 5y_2(t) + \frac{-20}{17}e^{-t} \\
(e^{-5t}y_2(t))' &= \frac{-20}{17}e^{-6t} \\
(e^{-5t}y_2(t)) &= \int \frac{-20}{17}e^{-6s}ds + c_2 \\
&= \frac{10}{51}e^{-6t}.
\end{aligned}
$$

Thus, isolating for $y_2(t)$, we find

$$y_2(t) = \frac{10}{51}e^{-t} + c_2 e^{5t}$$

Now, we substitute this into the top equation, getting that

$$y_1'(t) = 5y_1(t) + \frac{10}{51}e^{-t} + c_2 e^{5t} + \frac{21}{17}e^{-t}.$$

Using the same method of Integrating Factors, we derive the formula for $y_1(t)$.

$$
\begin{aligned}
y_1'(t) &= 5y_1(t) + \frac{10}{51}e^{-t} + c_2 e^{5t} + \frac{21}{17}e^{-t} \\
(e^{-5t}y_1(t))' &= \frac{73}{51}e^{-6t} + c_2 \\
(e^{-5t}y_1(t)) &= \int \frac{73}{51}e^{-6s} + c_2 ds + c_1 \\
&= \frac{73}{306}e^{-6t} + c_2 t + c_1
\end{aligned}
$$

Again, isolating for $y_1(t)$ gives the formula

$$y_1(t) = \frac{73}{306}e^{-t} + c_2 t e^{5t} + c_1 e^{5t}.$$

Changing bases back by $\vec{x} = S\vec{y}$, we arrive as a particular solution,

$$
\vec{x}(t) = \begin{pmatrix} 2 & 5/2 & -2 \\ 0 & 1/2 & -5/2 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{73}{306}e^{-t} + c_1 e^{5t} + c_2 t e^{5t} \\ \frac{10}{51}e^{-t} + c_2 e^{5t} \\ \frac{4}{17}e^{-t} + c_3 \end{pmatrix} = S\vec{y} + Se^{Jt}\vec{c}.
$$

Note that because we have kept the constants, $c_1, c_2, c_3$, this is also the general solution to the non-homogeneous equation. However, we could have let all those constants be 0. This would have simplified our calculations and we still would have gotten a particular solution to $\frac{d}{dt}\vec{x}(t) = A\vec{x} + \vec{g}(t)$. To get the general solution, we simply add the general solution to the corresponding homogeneous equation to our particular solution.

### Reflection Questions

1. What are the steps to the method of Integrating Factors that you learned in 307?

2. What are the steps to the method of integrating factors in the "Variation of Parameters" section, above?

3. In your own words, what are the steps for solving $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ by change of variables?

4. Could this process be simplified if $A$ were an upper triangular matrix? A lower triangular matrix?

5. Do solutions to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ form a vector space?

6. What does the general solution to $(\frac{d}{dt} - A)\vec{x}(t) = \vec{g}(t)$ look like? Can you prove this from the definition of a general solution?

# 3 Fourier Series

## 3.1 Basis coefficients and Orthogonality

Let's begin by reviewing some of the structures of $\mathbb{R}^n$. As we learned in 308, $\mathbb{R}^n$ has the structure of a vector space. That is,

1. For any $\vec{x}, \vec{y} \in \mathbb{R}^n$ and $c \in \mathbb{R}$, $\vec{x} + c\vec{y} \in \mathbb{R}^n$     (Closure under addition and scalar multiplication)

2. There exists a vector $\vec{0}$ such that $\vec{x} + \vec{0} = \vec{x}$ for all $\vec{x} \in \mathbb{R}^n$.

3. For all $\vec{x}$, there exists a vector, $-\vec{x}$ such that $\vec{x} + -\vec{x} = \vec{0}$.

4. The following algebraic rules hold for all $\vec{x}, \vec{y}, \vec{z} \in \mathbb{R}^n$ and scalars $a, b \in \mathbb{R}$:

   - $(\vec{x} + \vec{y}) = (\vec{y} + \vec{x})$
   - $(\vec{x} + \vec{y}) + \vec{z} = \vec{x} + (\vec{y} + \vec{z})$
   - $a(\vec{x} + \vec{y}) = a\vec{x} + a\vec{y}$
   - $a(b\vec{x}) = (ab)\vec{x}$
   - $(a + b)\vec{x} = a\vec{x} + b\vec{x}$

   For our purposes, **one of the most important properties as a vector space is that $\mathbb{R}^n$ has bases. For any basis, $\mathcal{B} = \{\vec{b}_1, \vec{b}_2, ..., \vec{b}_n\}$, we can write any vector, $\vec{x} \in \mathbb{R}^n$, as a unique linear combination**

$$\vec{x} = c_1\vec{b}_1 + c_2\vec{b}_2 + ....c_n\vec{b}_n.$$

   But, $\mathbb{R}^n$ has additional structure, not related to begin a vector space (that is, that don't have to do with scaling and adding vectors together). $\mathbb{R}^n$ **also has an inner product, also known as the dot product.**

**Definition 3.1.** Given a vector space, $V$, an **inner product** on $V$ is a function, $\langle \cdot, \cdot \rangle_V : V \times V \to \mathbb{R}$ obeying the following rules,

1. (Symmetry) For all $\vec{x}, \vec{y} \in V$,
$$\langle \vec{x}, \vec{y} \rangle_V = \langle \vec{y}, \vec{x} \rangle_V$$

2. (Bi-Linearity) For all $\vec{x}, \vec{y}, \vec{z} \in V$ and scalar $c \in \mathbb{R}$,
$$\langle \vec{x}, \vec{y} + c\vec{z} \rangle_V = \langle \vec{x}, \vec{y} \rangle_V + c\langle \vec{x}, \vec{z} \rangle_V$$

   and
$$\langle \vec{x} + c\vec{y}, \vec{z} \rangle_V = \langle \vec{x}, \vec{z} \rangle_V + c\langle \vec{y}, \vec{z} \rangle_V$$

3. (Positive-definite) For all $\vec{x} \in V$,
$$0 \leq \langle \vec{x}, \vec{x} \rangle_V$$
and $0 = \langle \vec{x}, \vec{x} \rangle_V$ if and only if $\vec{x} = \vec{0}$.

Recall that **an inner product lets us talk about angles.**

**Definition 3.2.** Two vectors, $\vec{x}, \vec{y} \in \mathbb{R}^n$ are called **orthogonal** in $\mathbb{R}^n$ if
$$\langle \vec{x}, \vec{y} \rangle_{\mathbb{R}^n} = 0.$$

Recall, we can define the norm or magnitude of a vector in terms of the inner product, by $\langle \vec{x}, \vec{x} \rangle_{\mathbb{R}^n} = |\vec{x}|^2$. In $\mathbb{R}^n$ this is just the Pythagorean Theorem.

**Definition 3.3.** Suppose that $\mathcal{B} = \{\vec{b}_1, \vec{b}_2, , ..., \vec{b}_n\}$ is an orthogonal basis for $\mathbb{R}^n$. Because it is a basis, we know that for every $\vec{x} \in \mathbb{R}^n$, there exists a unique expression in the basis, $\mathcal{B}$. But, how do we find the coefficients, $c_1, c_2, ..., c_n$ such that $\vec{x} = c_1, \vec{b}_1 + c_2 \vec{b}_2 + ... + c_n \vec{b}_n$?

**Answer:** We use the inner product. We know that there exist some constants such that
$$\vec{x} = c_1, \vec{b}_1 + c_2 \vec{b}_2 + ... + c_n \vec{b}_n$$

So, we can use the inner product and the fact that $\mathcal{B}$ is a collection of mutually orthogonal vectors to get

$$
\begin{aligned}
\langle \vec{x}, \vec{b}_1 \rangle_{\mathbb{R}^n} &= \langle c_1 \vec{b}_1 + c_2 \vec{b}_2 + ... + c_n \vec{b}_n, \vec{b}_1 \rangle_{\mathbb{R}^n} \\
\langle \vec{x}, \vec{b}_1 \rangle_{\mathbb{R}^n} &= c_1 \langle \vec{b}_1, \vec{b}_1 \rangle_{\mathbb{R}^n} + c_2 \langle \vec{b}_2, \vec{b}_1 \rangle_{\mathbb{R}^n} + ... + c_n \langle \vec{b}_n, \vec{b}_1 \rangle_{\mathbb{R}^n} \\
\langle \vec{x}, \vec{b}_1 \rangle_{\mathbb{R}^n} &= c_1 \langle \vec{b}_1, \vec{b}_1 \rangle_{\mathbb{R}^n}
\end{aligned}
$$

Thus we have that
$$c_1 = \frac{\langle \vec{x}, \vec{b}_1 \rangle_{\mathbb{R}^n}}{\langle \vec{b}_1, \vec{b}_1 \rangle_{\mathbb{R}^n}}.$$
and, in general,
$$c_i = \frac{\langle \vec{x}, \vec{b}_i \rangle_{\mathbb{R}^n}}{\langle \vec{b}_i, \vec{b}_i \rangle_{\mathbb{R}^n}}.$$

So, we now know that for any $\vec{x} \in \mathbb{R}^n$ and any orthogonal basis $\mathcal{B}$,

$$\vec{x} = \frac{\langle \vec{x}, \vec{b}_1 \rangle_{\mathbb{R}^n}}{\langle \vec{b}_1, \vec{b}_1 \rangle_{\mathbb{R}^n}} \vec{b}_1 + \frac{\langle \vec{x}, \vec{b}_2 \rangle_{\mathbb{R}^n}}{\langle \vec{b}_2, \vec{b}_2 \rangle_{\mathbb{R}^n}} \vec{b}_2 + ... + \frac{\langle \vec{x}, \vec{b}_n \rangle_{\mathbb{R}^n}}{\langle \vec{b}_n, \vec{b}_n \rangle_{\mathbb{R}^n}} \vec{b}_n.$$

**It is this structure which we wish to generalize to infinite dimensional vector spaces.**

## 3.2 Infinite-dimensional vector spaces

In this class, we will be interested in a very specific family of infinite-dimensional vector spaces.

**Definition 3.4.** For any interval, $[a, b] \in \mathbb{R}^n$ with $a < b$, we define the infinite-dimensional vector space $L^2([a, b], \mathbb{R})$ as follows:

$$L^2([a, b], \mathbb{R}) = \{f : [a, b] \to \mathbb{R} \quad | \quad \int_a^b |f|^2 dx < \infty\}.$$

Note that vectors in this vector space are functions defined on the interval $[a, b]$. Because of this, we will often refer to $L^2([a, b], \mathbb{R})$ as a **function space**, as well.

Whether a function is in $L^2([a, b], \mathbb{R})$ or not only depends upon what a function does on the interval $[a, b]$. For example, consider the function $f(x) = 1/x$. It is easy to check that $f \in L^2([1, 2], \mathbb{R})$, since

$$\int_1^2 \frac{1}{x^2} dx \leq 1$$

but $f \notin L^2([0, 2], \mathbb{R})$ since

$$\int_0^2 \frac{1}{x^2} dx = \frac{-1}{x}\Big|_0^2 = \infty.$$

**Question 3.2:** Is $L^2([a, b], \mathbb{R})$ really a vector space?

**Answer:** Yes. To see this, we need only check the definition of a vector space. We shall leave the details to the reader, but it is easy to see that the sum of two functions is still a function. Similarly with scalar multiplication. The zero vector is the zero function. And the additive inverse of a vector, $f$, is simply $-f$.

**Question 3.3:** Do infinite-dimensional vector spaces have bases? What do they look like?

**Answer:** Yes. All vector spaces have bases. However, for infinite-dimensional vector spaces, sometimes these bases are essentially useless. This stems from our definition of a basis. Recall that **a basis, $\mathcal{B}$, is a linearly independent collection of vectors such that every vector can be expressed (uniquely) as a finite linear combination of the vectors in $\mathcal{B}$.** So, consider the following example.

**Example 3.1.** Consider $\ell^1(\mathbb{R}) = \{(x_1, x_2, ...) : \sum_{i=1}^{\infty} x_i < \infty\}$ the space of infinite stings of real numbers which are square summable. It is left as an exercise to check that this is, in fact, a vector space.

Based upon our intuitions of finite-dimensional vector spaces like $\mathbb{R}^n$, we might guess that if we let

$$
\begin{aligned}
\vec{e}_1 &= (1, 0, 0, ....) \\
\vec{e}_2 &= (0, 1, 0, ....) \\
&\vdots \\
\vec{e}_n &= (0, 0, ...., 0, 1, 0, ...)
\end{aligned}
$$

then, the collection $\mathcal{B} = \{\vec{e}_1, \vec{e}_2, \vec{e}_3, ...\}$ might be a basis. But it isn't. The only vectors which are finite linear combinations of $\mathcal{B}$ are ones which have only finitely many non-zero entries. But the vector

$$\vec{x} = (1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, ...) \in \ell^1(\mathbb{R}).$$

Since $\vec{x}$ has infinitely many non-zero entries, it is not a finite linear combination of $\mathcal{B}$. There is a basis for $\ell^1(\mathbb{R})$, but it is essentially useless. Instead, we will use a different idea.

**Definition 3.5.** For a vector space, $V$, a **countable basis** is a collection of linearly independent vectors, $\mathcal{B}$ such that every vector, $\vec{v} \in V$, can be arbitrarily well-approximated by finite linear combinations of $\mathcal{B}$. More precisely, a countable basis is a collection of linearly independent vectors, $\mathcal{B} = \{\vec{b}_1, \vec{b}_2, \vec{b}_3, ...\}$ such that for every vector $\vec{v} \in V$ and every $\epsilon > 0$, there exists a finite linear combination of vectors in $\mathcal{B}$ such that

$$|\vec{v} - \sum_{i=1}^{N} c_i \vec{b}_i| < \epsilon.$$

Note that in finite-dimensional vector spaces like $\mathbb{R}^n$, a basis and a countable basis are the same thing. However, in infinite-dimensional vector spaces, they are very different. Let's us return to our example.

**Example 3.2.** If we let

$$
\begin{aligned}
\vec{e}_1 &= (1, 0, 0, 0, ....) \\
\vec{e}_2 &= (0, 1, 0, 0, ....) \\
\vec{e}_3 &= (0, 0, 1, 0....) \\
&\vdots
\end{aligned}
$$

then, the collection $\mathcal{B} = \{\vec{e}_1, \vec{e}_2, \vec{e}_3, ...\}$ is a countable basis for the vector space $\ell^1(\mathbb{R}) = \{(x_1, x_2, ...) : \sum_{i=1}^{\infty} x_i < \infty\}$.

To see this, we remark that for any $\vec{v} \in \ell^1(\mathbb{R})$ and any $\epsilon > 0$ we can find an integer such that if we write $\vec{v} = (v_1, v_2, ....)$, then

$$\sum_{i=N}^{\infty} v_i < \epsilon.$$

Thus, $|\sum_{i=1}^{N} v_i \vec{e}_i - \vec{v}| < \epsilon$.

In this class, we will use countable bases to get a hold of infinite-dimensional vector spaces.

**Question 3.4:** What is a countable basis for $L^2([a,b], \mathbb{R})$?

**Answer:** As with all bases, there are infinitely many. However, in this class, we will only deal with one, special countable basis. The collection of functions,

$$\{cos(n\frac{2\pi}{b-a}x)\}_{n=0}^{\infty} \cup \{sin(n\frac{2\pi}{b-a}x)\}_{n=1}^{\infty}$$

is a countable basis for the function space $L^2([a,b], \mathbb{R})$. We can think of it as a sort of standard basis for $L^2([a,b], \mathbb{R})$.

This means that every function in $L^2([a,b], \mathbb{R})$ can be arbitrarily well-approximated by finite linear combinations of sine and cosine functions. And, **in some sense, for every** $f \in L^2([a,b], \mathbb{R})$ **there exists an infinite linear combination such that**

$$f(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x)$$

**Definition 3.6.** The basis expression of a function, $f \in L^2([a,b], \mathbb{R})$ in terms of the basis $\{cos(n\frac{2\pi}{b-a}x)\}_{n=0}^{\infty} \cup \{sin(n\frac{2\pi}{b-a}x)\}_{n=1}^{\infty}$ is called the **Fourier Series of** $f$ **on the interval** $[a,b]$. We will write it as

$$\mathcal{F}_{[a,b]}(f)(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x)$$

**Question 3.5:** Can we define an inner product on $L^2([a,b], \mathbb{R})$?

**Answer:** Yes. Having an inner product is a special structure that cannot always be put on infinite-dimensional spaces, but the spaces $L^2([a,b], \mathbb{R})$ do have an inner product.

**Definition 3.7.** We define the **inner product on** $L^2([a,b], \mathbb{R})$ as follows. For any two functions (vectors), $f, g \in L^2([a,b], \mathbb{R})$,

$$\langle f, g \rangle_{L^2([a,b], \mathbb{R})} = \int_a^b f(x)g(x)dx.$$

Recall that with an inner product, we can define the **norm** or magnitude of a vector. Thus, for a function (vector), $f$, in $L^2([a, b], \mathbb{R})$, we define it's norm by

$$||f||_{L^2([a,b],\mathbb{R})} = \langle f, f \rangle_{L^2([a,b],\mathbb{R})}^{1/2} = (\int_a^b |f|^2 dx)^{1/2}$$

Thus, the condition in the definition of $L^2([a, b], \mathbb{R})$ about $\int_a^b |f|^2 dx < \infty$ is really a condition that the norm of vectors be finite. That is, we only want to consider vectors for which $||f||_{L^2([a,b],\mathbb{R})} < \infty$

**Question 3.6:** Is the basis $\{cos(n\frac{2\pi}{b-a}x)\}_{n=0}^{\infty} \cup \{sin(n\frac{2\pi}{b-a}x)\}_{n=1}^{\infty}$ an orthogonal basis in $L^2([a, b], \mathbb{R})$?

**Answer:** Yes. This is left as a homework exercise.

**Question 3.7:** Given a function, $f : [a, b] \to \mathbb{R}$, in $L^2([a, b], \mathbb{R})$ how do we find the basis coefficients, $\{a_n\}_{n=0}^{\infty} \cup \{b_n\}_{n=1}^{\infty}$ of $\mathcal{F}_{[a,b]}(f)$?

**Answer:** As always, given an orthogonal basis, we find the basis coefficients by taking the inner product with the basis element in question.

$$
\begin{aligned}
\mathcal{F}_{[a,b]}(f)(x) &= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x) \\
\langle \mathcal{F}_{[a,b]}(f)(x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} &= \langle \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} \\
\langle \mathcal{F}_{[a,b]}(f)(x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} &= \frac{a_0}{2} \langle 1, cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} \\
&\quad + \sum_{n=1}^{\infty} a_n \langle cos(n\frac{2\pi}{b-a}x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} \\
&\quad + \sum_{n=1}^{\infty} b_n \langle sin(n\frac{2\pi}{b-a}x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} \\
\langle \mathcal{F}_{[a,b]}(f)(x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})} &= a_n \langle cos(n\frac{2\pi}{b-a}x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})}
\end{aligned}
$$

Thus,

$$a_n = \frac{\langle \mathcal{F}_{[a,b]}(f)(x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})}}{\langle cos(n\frac{2\pi}{b-a}x), cos(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})}}$$

Recalling that $f = \mathcal{F}_{[a,b]}(f)$ in the vector space $L^2([a, b], \mathbb{R})$ and the definition of the inner product, we have the equivalent definition,

$$a_n = \frac{\int_a^b f(x)cos(n\frac{2\pi}{b-a}x)dx}{\int_a^b |cos(n\frac{2\pi}{b-a}x)|^2 dx}$$

Similarly, we have the formulae for the $b_n$,

$$b_n = \frac{\langle \mathcal{F}_{[a,b]}(f)(x), sin(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})}}{\langle sin(n\frac{2\pi}{b-a}x), sin(n\frac{2\pi}{b-a}x) \rangle_{L^2([a,b],\mathbb{R})}}$$

or

$$b_n = \frac{\int_a^b f(x)sin(n\frac{2\pi}{b-a}x)dx}{\int_a^b |sin(n\frac{2\pi}{b-a}x)|^2dx}.$$

**Question 3.8:** The denominators, $\int_a^b |sin(n\frac{2\pi}{b-a}x)|^2dx$ and $\int_a^b |cos(n\frac{2\pi}{b-a}x)|^2dx$, seem complicated and ugly. What are these?

**Answer:** Fortunately, for $n \neq 0$,

$$\int_a^b |sin(n\frac{2\pi}{b-a}x)|^2dx = \int_a^b |cos(n\frac{2\pi}{b-a}x)|^2dx = \frac{b-a}{2}.$$

When $n = 0$, $\int_a^b 1dx = b - a$.

Thus, our formulae for the basis coefficients of $\mathcal{F}_{[a,b]}(f)$ become, for all $a_n$ and $b_n$

$$a_n = \frac{2}{b-a}\int_a^b f(x)cos(n\frac{2\pi}{b-a}x)dx$$

$$b_n = \frac{2}{b-a}\int_a^b f(x)sin(n\frac{2\pi}{b-a}x)dx$$

These are called the **Fourier-Euler equations.**

**Example 3.3.** Let $f : [0, 2\pi] \to \mathbb{R}$ be defined piecewise by

$$f(x) = \begin{cases} 1 & x \in [0, \pi] \\ 0 & x \in (\pi, 2\pi] \end{cases}$$

Find $\mathcal{F}_{[0,2\pi]}(f)(x)$, explicitly.

First, is $f \in L^2([0, 2\pi], \mathbb{R})$? We check from the definition:

$$\int_0^{2\pi} |f|^2dx < \int_0^{2\pi} 1dx \leq 2\pi < \infty.$$

So, $f \in L^2([0, 2\pi], \mathbb{R})$. Thus, we know that $\{cos(nx)\}_{n=0}^\infty \cup \{sin(nx)\}_{n=1}^\infty$ is an orthogonal, countable basis for $L^2([0, 2\pi], \mathbb{R})$. So, we can find a basis expression for $f$

$$\mathcal{F}_{[0,2\pi]}(f)(x) = \frac{a_0}{2} + \sum_{n=1}^\infty a_n cos(nx) + b_n sin(nx)$$

where the coefficients are given by the Fourier-Euler formulae

$$a_n = \frac{1}{\pi}\int_0^{2\pi} f(x)cos(nx)dx$$

61

$$b_n = \frac{1}{\pi} \int_0^{2\pi} f(x) sin(nx) dx.$$

Thus, we need only calculate the coefficients. We begin with $a_n$.

$$
\begin{aligned}
a_n &= \frac{1}{\pi} \int_0^{2\pi} f(x) cos(nx) dx \\
&= \frac{1}{\pi} [\int_0^{\pi} f(x) cos(nx) dx + \int_{\pi}^{2\pi} f(x) cos(nx) dx] \\
&= \frac{1}{\pi} \int_0^{\pi} cos(nx) dx \\
&= \frac{1}{\pi n} sin(nx)|_0^{\pi} \\
&= 0
\end{aligned}
$$

Note that this integration only works for $n > 0$, so we must calculate $a_0$ separately.

$$
\begin{aligned}
a_0 &= \frac{1}{\pi} \int_0^{2\pi} f(x) cos(0x) dx \\
&= \frac{1}{\pi} \int_0^{2\pi} f(x) dx \\
&= \frac{1}{\pi} \int_0^{\pi} 1 dx \\
&= 1
\end{aligned}
$$

Now we calculate the $b_n$

$$
\begin{aligned}
b_n &= \frac{1}{\pi} \int_0^{2\pi} f(x) sin(nx) dx \\
&= \frac{1}{\pi} [\int_0^{\pi} f(x) sin(nx) dx + \int_{\pi}^{2\pi} f(x) sin(nx) dx] \\
&= \frac{1}{\pi} \int_0^{\pi} sin(nx) dx \\
&= \frac{-1}{\pi n} cos(nx)|_0^{\pi} \\
&= \frac{-1}{\pi n} [1 - (-1)^n]
\end{aligned}
$$

Plugging into our formula for the Fourier Series,

$$\mathcal{F}_{[0,2\pi]}(f)(x) = \frac{1}{2} + \sum_{n=1}^{\infty} \frac{-1}{\pi n} [1 - (-1)^n] sin(nx).$$

## 3.3 What does $\mathcal{F}_{[a,b]}(f)(x)$ look like?

**Question 3.9:** Earlier, I said that "in some sense" $f = \mathcal{F}_{[a,b]}(f)$. In what sense?

**Answer:** There are two senses in which this this equality is funny. First, **it is not true that $f(x) = \mathcal{F}_{[a,b]}(f)(x)$ for all $x \in [a,b]$. But it is true that $f = \mathcal{F}_{[a,b]}(f)$ as vectors in the vector space $L^2([a,b], \mathbb{R})$**. That is, $||f - \mathcal{F}_{[a,b]}(f)||_{L^2([a,b],\mathbb{R})} = 0$. What does this mean? It means that

$$\int_a^b |f - \mathcal{F}_{[a,b]}(f)|^2 dx = 0.$$

So, if we let $G = \{x \in [a,b] : f(x) = \mathcal{F}_{[a,b]}(f)(x)\}$ and $B = \{x \in [a,b] : f(x) \neq \mathcal{F}_{[a,b]}(f)(x)\}$, then we can see that

$$
\begin{aligned}
0 &= \int_a^b |f - \mathcal{F}_{[a,b]}(f)|^2 dx \\
&= \int_G |f - \mathcal{F}_{[a,b]}(f)|^2 dx + \int_B |f - \mathcal{F}_{[a,b]}(f)|^2 dx \\
&= \int_B |f - \mathcal{F}_{[a,b]}(f)|^2 dx
\end{aligned}
$$

Since we have no control over the difference, $|f(x) - \mathcal{F}_{[a,b]}(f)(x)|$, all that we can say is that the set, $B$, must be so small that integrals cannot see it. Later, we will be interested in conditions on the function, $f$, such that we can control the bad set, $B$. But, for now, all we can say is that it is invisible to integrals.

The other sense in which the equation, $f = \mathcal{F}_{[a,b]}(f)$, is strange is that since it is only true in $L^2([a,b], \mathbb{R})$, it does not hold away from the interval $[a, b]$. For example, $f$ may have not been defined off the interval, $[a, b]$. On the other hand, no matter what $f$ is away from $[a, b]$, $\mathcal{F}_{[a,b]}(f)$ is defined on all of $\mathbb{R}$.

**Question 3.10:** What does $\mathcal{F}_{[a,b]}(f)$ look like outside of $[a, b]$?

**Answer:** Let's write out $\mathcal{F}_{[a,b]}(f)(x)$ from the definition.

$$
\mathcal{F}_{[a,b]}(f)(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x)
$$

So, for $x + (b - a)$,

$$
\begin{aligned}
\mathcal{F}_{[a,b]}(f)(x + (b-a)) &= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}(x + (b-a))) + b_n sin(n\frac{2\pi}{b-a}(x + (b-a))) \\
&= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x + 2n\pi)) + b_n sin(n\frac{2\pi}{b-a}x + 2n\pi) \\
&= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{2\pi}{b-a}x) + b_n sin(n\frac{2\pi}{b-a}x) \\
&= \mathcal{F}_{[a,b]}(f)(x)
\end{aligned}
$$

Thus, $\mathcal{F}_{[a,b]}(f)(x) = \mathcal{F}_{[a,b]}(f)(x + (b-a))$ **for all** $x \in \mathbb{R}$.

**Definition 3.8.** A function, $f$, is $P-$**periodic**, if $f(x) = f(x + P)$ for all $x \in \mathbb{R}$, and $P > 0$ is the smallest such number for which is this true.

In this class, for ease we will ignore this last part, that $P$ must be the smallest such positive number. Thus we will say that $\mathcal{F}_{[a,b]}(f)(x)$ **is** $(b - a)$**-periodic.**

**Question 3.11:** Under what conditions can we say concretely that $f(x) = \mathcal{F}_{[a,b]}(f)(x)$?

**Answer:** In general, we have no idea for which $x \in [a, b]$ $f(x) = \mathcal{F}_{[a,b]}(f)(x)$. In order to say anything, we must make much stronger assumptions on the function $f$. In particular, we will need both $f$ and $f'$ to be **piecewise continuous**.

**Definition 3.9.** A function, $f : [a, b] \to \mathbb{R}$, is called **piecewise continuous on** $[a, b]$ if we can cut up the interval $[a, b]$, into finitely many non-degenerate subintervals, $[x_i, x_{i+1}]$, where

$$a = x_0 < x_1 < ... < x_{n-1} < x_n = b$$

such that $f$ is continuous on each subinterval, $(x_i, x_{i+1})$ and for each endpoint of each subinterval, $lim_{y \to x_i^{\pm}} f(y) < \infty$.

Note that piecewise continuous functions can only have finitely many jump discontinuities. A removable discontinuity violates the assumption that the partition of $[a, b]$ must be non-degenerate (i.e., that $x_i < x_{i+1}$). Vertical asymptotes violate the assumption that $lim_{y \to x_i^{\pm}} f(y) < \infty$.

**Theorem 3.1.** *If a function $f \in L^2([a, b], \mathbb{R})$ is piecewise continuous and its derivative, $f'$ is also piecewise continuous, then*

$$\mathcal{F}_{[a,b]}(f)(x) = \frac{1}{2}[lim_{y \to x^+} f(y) + lim_{y \to x^-} f(y)]$$

Note that if $f$ is continuous at $x \in [a, b]$, then these limits are the same and both equal $f(x)$. Thus, $f(x) = \mathcal{F}_{[a,b]}(f)(x)$.

**Example:** Consider the function $f(x) = x^2$.

1. Describe $\mathcal{F}_{[0,1]}(f)(x)$.

   Since $\mathcal{F}_{[0,1]}(f)(x)$ converges to $f(x) = x^2$ in the interval $[0, 1]$, $f$ and $f'$ are piece-wise continuous on $[0, 1]$, and $\mathcal{F}_{[0,1]}(f)(x)$ is 1-periodic, we know that

   $$\mathcal{F}_{[0,1]}(f)(x) = \begin{cases} x^2 & \text{on } (0, 1) \\ \frac{1}{2} & x = 1 \end{cases} \qquad \mathcal{F}_{[0,1]}(f)(x) = \mathcal{F}_{[0,1]}(f)(x + 1)$$

   where $\mathcal{F}_{[0,1]}(f)(1) = \frac{1}{2}$ because of Theorem 3.1. Try drawing this function yourself.

2. Describe $\mathcal{F}_{[-1,1]}(f)(x)$.

   By the same reasoning, we have that

   $$\mathcal{F}_{[-1,1]}(f)(x) = x^2 \text{ on } [-1, 1], \qquad \mathcal{F}_{[-1,1]}(f)(x) = \mathcal{F}_{[-1,1]}(f)(x + 2)$$

   Try drawing this function yourself.

3. Describe $\mathcal{F}_{[-1,3]}(f)(x)$.

Again, because $\mathcal{F}_{[-1,3]}(f)(x)$ converges to $f(x) = x^2$ in the interval $[-1, 3]$, $f$ and $f'$ are piece-wise continuous on $[-1, 3]$, and $\mathcal{F}_{[-1,3]}(f)(x)$ is 4-periodic, we know that

$$\mathcal{F}_{[-1,3]}(f)(x) = \begin{cases} x^2 & \text{on } (-1,3) \\ \frac{10}{2} & x = 3 \end{cases} \qquad \mathcal{F}_{[-1,3]}(f)(x) = \mathcal{F}_{[-1,3]}(f)(x+4)$$

where $\mathcal{F}_{[0,1]}(f)(1) = \frac{10}{2}$ because of Theorem 3.1.

Try drawing this function yourself.

## 3.4    Even and Odd functions and Extensions

**Note that in the examples above, where $f(x) = x^2$, we saw that $\mathcal{F}_{[-1,3]}(f)(x) = \mathcal{F}_{[-1,1]}(f)(x) = \mathcal{F}_{[0,1]}(f)(x)$ for all points $x \in (0,1)$. These three different Fourier series, which all converge to different functions on $\mathbb{R}$, all converge to the same function on $(0,1)$.**

**Definition 3.10.** Given a function, $f : [a, b] \to \mathbb{R}$, we say that another function, $\widetilde{f} : [c, d] \to \mathbb{R}$, is an **extension of $f$ to the interval $[c, d]$** if

- The interval $[a, b]$ is contained in the interval $[c, d]$. That is, $c \le a$ and $b \le d$.

- For all points, $x \in [a, b]$, the functions agree, $f(x) = \widetilde{f}(x)$

What we saw in the above examples is that $\mathcal{F}_{[-1,3]}(f)(x), \mathcal{F}_{[-1,1]}(f)(x)$, and $\mathcal{F}_{[0,1]}(f)(x)$ are all extensions of the function $f(x) = x^2$ restricted to the interval $(0, 1)$.

**Question 3.12:** How are extensions of a function useful for Fourier Series?

**Answer:** Let's suppose we have a function, $f : [0, L] \to \mathbb{R}$, in $L^2([0, L], \mathbb{R})$ such that $f, f'$ are piecewise continuous. We could choose to find the Fourier Series on $[0, L]$, $\mathcal{F}_{[0,L]}(f)$. Or, we could choose to extend the function, $f$, to a new function, $\widetilde{f}$, defined on a larger interval, say, $[-L, L]$. Then we know that $\mathcal{F}_{[0,L]}(f)$ converges to $f$ and $\mathcal{F}_{[-L,L]}(\widetilde{f})$ converges to $\widetilde{f}$. Since $f$ and $\widetilde{f}$ are the same function on $[0, L]$, then, we have that

$$\mathcal{F}_{[0,L]}(f)(x) = \mathcal{F}_{[-L,L]}(\widetilde{f})(x) \quad \forall x \in (0, L)$$

Thus, we could choose ANY extension to ANY larger interval, and take the Fourier Series on that interval. This Fourier Series would converge on the smaller interval to the function

we extended. **It turns out, by choosing a clever extension, we can make some coefficients zero.**

Recall the following definitions.

**Definition 3.11.** A function, $f : [-a, a] \to \mathbb{R}$, defined on a **symmetric interval** is called even if

$$f(-x) = f(x)$$

and it is called **odd** if

$$f(-x) = -f(x)$$

for all $x \in [-a, a]$.

Recall also that

- Products of even functions are even.

- The product of two odd functions is an even function.

- The product of an off function with an even function is an odd function.

- If $f$ is odd, then

$$\int_{-a}^{a} f(x)dx = 0$$

- If $f$ is even

$$\int_{-a}^{a} f(x)dx = 2\int_{0}^{a} f(x)dx$$

**Question 3.12:** So, what happens if we choose an odd extension? That is, if $f : [0, L] \to \mathbb{R}$ and we define an odd function, $f^{odd} : [-L, L] \to \mathbb{R}$ by

$$f^{odd}(x) = \begin{cases} f(x) & [0, L] \\ -f(-x) & [-L, 0] \end{cases}$$

what does $\mathcal{F}_{[-L,L]}(f^{odd})(x)$ look like?

**Answer:** Well, to begin with, we use the formulas from the definition.

$$\mathcal{F}_{[-L,L]}(f^{odd})(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{\pi}{L}x) + b_n sin(n\frac{\pi}{L}x)$$

where the coefficients are given by the formulae

$$a_n = \frac{1}{L}\int_{-L}^{L} f^{odd}(x)cos(n\frac{\pi}{L}x)dx$$

66

$$b_n = \frac{1}{L} \int_{-L}^{L} f^{odd}(x)sin(n\frac{\pi}{L}x)dx.$$

But, since *sine* is an odd function and *cosine* is an even function, $f^{odd}(x)cos(nx)$ is an odd function and $f^{odd}(x)sin(nx)$ is an even function. Thus, $a_n = 0$ for all $n = 0, 1, 2, ....$ This means we only need to compute the $b_n$. Even better, because $f^{odd}(x)sin(n\frac{\pi}{L}x)$ is an even function, though, we can simplify the formula to

$$b_n = \frac{2}{L} \int_{0}^{L} f(x)sin(n\frac{\pi}{L}x)dx.$$

$\mathcal{F}_{[-L,L]}(f^{odd})(x)$ is called the **Fourier Sine Series of** $f$.

**Question 3.12:** So, what happens if we choose an even extension? That is, if $f : [0, L] \rightarrow \mathbb{R}$ and we define an even function, $f^{even} : [-L, L] \rightarrow \mathbb{R}$ by

$$f^{even}(x) = \begin{cases} f(x) & [0, L] \\ f(-x) & [-L, 0] \end{cases}$$

what does $\mathcal{F}_{[-L,L]}(f^{even})(x)$ look like?

**Answer:** Again, we use the formulas from the definition.

$$\mathcal{F}_{[-L,L]}(f^{even})(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{\pi}{L}x) + b_n sin(n\frac{\pi}{L}x)$$

where the coefficients are given by the formulae

$$a_n = \frac{1}{L} \int_{-L}^{L} f^{even}(x)cos(n\frac{\pi}{L}x)dx$$

$$b_n = \frac{1}{L} \int_{-L}^{L} f^{even}(x)sin(n\frac{\pi}{L}x)dx.$$

But now, since $f^{even}(x)sin(n\frac{\pi}{L}x)$ is odd, all of the $b_n = 0$. And, since $f^{even}(x)cos(n\frac{\pi}{L}x)$ are even, the formulae for the $a_n$ can be simplified to

$$a_n = \frac{2}{L} \int_{0}^{L} f(x)cos(n\frac{\pi}{L}x)dx.$$

$\mathcal{F}_{[-L,L]}(f^{even})(x)$ is called the **Fourier Cosine Series of** $f$.

Thus, in choosing a smart extension, we can simplified our work and only compute half of the coefficients.

**Example 3.4.** Let $f = \begin{cases} L & x \in [0, L] \\ x & x \in (L, 2L] \end{cases}$. Give a Fourier Sine Series for $f$.

First, since we want a Sine Series, so we must take the odd extension of $f$. Then, we have

$$\mathcal{F}_{[-2L,2L]}(f^{odd})(x) = \sum_{n=1}^{\infty} b_n sin(n\frac{\pi}{2L}x)$$

where

$$b_n = \frac{1}{L} \int_0^{2L} f(x)sin(n\frac{\pi}{2L}x)dx.$$

Now, we compute the coefficients.

$$
\begin{aligned}
b_n &= \frac{1}{L} \int_0^{2L} f(x)sin(n\frac{\pi}{2L}x)dx \\
&= \frac{1}{L} \int_0^L Lsin(n\frac{\pi}{2L}x)dx + \frac{1}{L} \int_L^{2L} xsin(n\frac{\pi}{2L}x)dx \\
&= \frac{-2}{n\pi}cos(n\frac{\pi}{2L}x)|_0^L + \frac{-2}{n\pi}xcos(n\frac{\pi}{2L}x)|_L^{2L} + \frac{-2}{n\pi} \int_L^{2L} cos(n\frac{\pi}{2L}x)dx \\
&= \frac{-2}{n\pi}[cos(n\frac{\pi}{2}) - 1] + \frac{-2}{n\pi}[2L(-1)^n - Lcos(n\frac{\pi}{2})] \\
&\quad + \frac{-L}{n\pi^2}[-sin(n\frac{\pi}{2})]
\end{aligned}
$$

Plugging in, we have

$$\mathcal{F}_{[-2L,2L]}(f^{odd})(x) = \sum_{n=1}^{\infty} [\frac{-2}{n\pi}[cos(n\frac{\pi}{2}) - 1] + \frac{-2}{n\pi}[2L(-1)^n - Lcos(n\frac{\pi}{2})] + \frac{L}{n\pi^2}sin(n\frac{\pi}{2})]sin(n\frac{\pi}{2L}x)$$

### Reflection Questions

1. What does it mean for a function, $g$, to be an extension of a function, $f$?

2. What is the definition of an even function? of an odd function?

3. Let $f : [0, L] \to \mathbb{R}$ be given by $f(x) = x$.

   a. Is $f$ an even function, an odd function, or neither?

   b. Draw two different even extensions of $f$.

   c. Draw two different odd extensions of $f$.

   d. In your own words, what steps would you take to find the Fourier Series of each of the extensions of $f$, above?

   e. Use Theorem 3.1, above, to draw the functions that those Fourier Series would converge to.

4. Show that the product of an even function and an odd function is an odd function.

5. Is $tan(x)$ a piece-wise continuous function?

6. For any vector space, there are infinitely many bases. If we can find Fourier Sine series and Fourier Cosine series for any function $f \in L^2([0, L], \mathbb{R})$, does that mean that both

$$\{sin(n\frac{\pi}{L}x)\}_{n=1}^{\infty} \quad \text{and} \quad \{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$$

are each countable bases for $L^2([0, L], \mathbb{R})$? How does this differ from the standard orthogonal countable basis we have been using?

# 4 Boundary Value Problems

**Example 4.1..** Consider the differential equation

$$y'' + \lambda y = 0$$

on the domain $[0, L] \subset \mathbb{R}$. What if we wanted to find all the solutions to this equation which also satisfied the boundary conditions $y(0) = 0$ and $y(L) = 0$ for some fixed $L > 0$?

This looks similar to initial value problems, so let's recall how we solve those: first, we find the general solution to the differential equation, then we plug in to determine the coefficients. Let's try it.

First, we find the general solution to the differential equation. Note that this depends upon the value of $\lambda$.

1. Case 1: $\lambda < 0$. Then the general solution is

$$y(t) = c_1 e^{\sqrt{-\lambda}t} + c_2 e^{-\sqrt{-\lambda}t}$$

2. Case 2: $\lambda = 0$. In this case the general solution is

$$y(t) = c_1 t + c_2$$

3. Case 3: $\lambda > 0$. In this case, the general solution is

$$y(t) = c_1 cos(\sqrt{\lambda}t) + c_2 sin(\sqrt{\lambda}t)$$

Now, we compare our general solution with our boundary conditions. That is, we search for constant $c_1, c_2$ such that $y(t)$ satisfies $y(0) = 0$ and $y(L) = 0$.

1. Case 1: $\lambda < 0$. Checking the first endpoint, we see that

$$y(0) = c_1 + c_2 = 0$$

Thus, $c_1 = -c_2$. Checking the other endpoint, we have

$$y(L) = c_1(e^{\sqrt{-\lambda}L} - e^{-\sqrt{-\lambda}L}) = 0.$$

Since $L > 0$, $e^{\sqrt{-\lambda}L} > e^{-\sqrt{-\lambda}L}$, and thus $c_1 = 0$.

This means that for $\lambda < 0$, there are NO non-trivial solutions to this problem! Only the zero function solves the problem.

2. Case 2: $\lambda = 0$. Checking the first endpoint,

$$y(0) = c_2 = 0$$

Checking the second endpoint, we see that

$$y(L) = c_1 L = 0.$$

Therefore for $\lambda = 0$, the only solution is the trivial solution.

3. Case 3: $\lambda > 0$. Checking the first endpoint,

$$y(0) = c_1 cos(\sqrt{\lambda}0) + c_2 sin(\sqrt{\lambda}0) = c_1 = 0$$

Checking the second endpoint, we see that

$$y(L) = c_2 sin(\sqrt{\lambda}L) = 0$$

This means that either $c_2 = 0$ or $sin(\sqrt{\lambda}L) = 0$. If $c_2 = 0$, then $x(t)$ is the zero function. This is uninteresting. Notice that since the differential equation is homogeneous and the boundary conditions are all zero, we always have the zero function as a solution. On the other hand, $sin(\sqrt{\lambda}L) = 0$ only if $\sqrt{\lambda}L = n\pi$ for some integer $n$. That is, if $\lambda = \frac{n^2\pi^2}{L^2}$, then

$$y(t) = c_2 sin(\frac{n\pi}{L}t)$$

is a solution for any value of $c_2$. **This means that we would have infinitely many solutions.**

   This is very different from initial value problems where we had one solution. These types of problems are called Boundary Value Problems.

**Definition 4.1. A boundary value problem (BVP)** consists of three things:

1. a differential equation,

2. a domain, $\Omega$,

3. and boundary conditions on the boundary of $\Omega$.

**Definition 4.2. A solution to a boundary value problem** is a function, $f : \Omega \to \mathbb{R}$, such that

1. $f$ is a solution to the differential equation in $\Omega$,

2. $f$ agrees with the boundary conditions on the boundary of $\Omega$.

   Returning to the previous example, observe that we can think of this problem as looking for values of $\lambda$ for which there are non-zero functions $y$ which are scaled by the action of $\frac{d^2}{dt^2}$,

$$\frac{d^2}{dt^2}y = -\lambda y$$

and satisfy the boundary conditions. In some sense, **this is an eigenvalue problem.**

**Definition 4.3.** For any linear operator, $L$, a non-trivial function, $y$, which satisfies the equation $L(y) = \lambda y$ is called an **eigenfunction** of the linear operator $L$. If such an eigenfunction exists, $\lambda$ is called an **eigenvalue** of the linear operator $L$.

In the example above, $\frac{d^2}{dt^2}$ is our linear operator. This means that if $y'' + \lambda y = 0$ for a non-zero function $y$, then $-\lambda$ is an eigenvalue of the linear operator $\frac{d^2}{dt^2}$. Since we can solve $y'' + \lambda y = 0$ for any $\lambda$, the set of all eigenvalues of $\frac{d^2}{dt^2}$ is all of $\mathbb{R}$. This is too large to deal with for us in this class. So, we impose our boundary conditions.

**Definition 4.4.** Let $L$ be a linear operator. A non-trivial function, $y$, which satisfies the equation $L(y) = \lambda y$ AND satisfies the boundary conditions of a BVP is called an **eigenfunction** of the BVP. If such an eigenfunction exists, $\lambda$ is called an **eigenvalue** of the BVP.

In the example above, the only eigenvalues of the BVP were $\lambda = -\frac{n^2\pi^2}{L^2}$ and the only eigenfunctions were of the form $y(t) = c_2 sin(\frac{n\pi}{L}t)$.

It is absolutely possible to ask if there is a "best basis" by which to view the action of the linear operator $\frac{d^2}{dt^2}$ and find the (infinite) matrix which represents the action of the linear operator $\frac{d^2}{dt^2}$ on the infinite dimensional vector space of twice-differentiable functions on $[0, L]$. However, we will not make this approach explicit in this class.

We now return to investigating two-point BVP.

**Example 4.2** A BVP need not have any solutions at all. If we consider the following boundary value problem,
$$y'' = 0$$
on the domain $[0, L] \subset \mathbb{R}$ with the boundary conditions $y'(0) = 0$ and $y'(L) = 2$.

The general solution to the differential equation is $x(t) = c_1 t + c_2$. But, affine linear functions have constant derivatives, so no affine linear function can satisfy the boundary conditions.

**Question 4.1:** Do the set of solutions to boundary value problems have any structure?

**Answer:** Yes. But, it is not as clean as we saw earlier for Initial Value Problems. First, a few definitions.

**Definition 4.5.** A boundary value problem is called **homogeneous** if the boundary conditions are zero and the differential equation is homogeneous. A boundary value problem is called **non-homogenous** otherwise.

**As usual, the structure of the set of solutions in the non-homogeneous case depends upon the structure of the set of solutions for the corresponding homogenous case.** That is, a non-homogeneous boundary value problem may have no, a unique, or infinitely many solutions. And, a homogeneous BVP always has at least the trivial solution, but may have infinitely many solutions.

**Theorem 4.1.** *(Structure of the set of solutions for BVP) A non-homogeneous BVP has a unique solution if and only if the corresponding homogeneous BVP has only the trivial solution. It has either no or infinitely many solutions if and only if the corresponding BVP has infinitely many solutions.*

## 4.1 Boundary Value Problems for the Heat Equation

The heat equation is of profound importance to the development of mathematics, physics, and still drives a lot of research in modern math. We will study it in the following context:

**Suppose that I made you this bet: I will take a metal wire of length, $L$, and heat it up in some places so that there is some initial heat distribution on it. If you can tell me what the heat distribution will be 10 minutes later, based only on the initial heat distribution, then I will give you a 4.0 in the class.**

**Question 4.2:** How would you model this problem using mathematics?

**Answer:** There are many answers to this question. But let's begin with modeling the wire. Let $x$ represent the spatial dimension. Thus, we can think of the points $x \in [0, L]$ as representing points on our wire.

Now, let's think of what we want: we want a function, $u(x, t)$ such that $u(x, t)$ tells us how much heat is at the point $x$ on our wire at time $t$. Thus, the domain of our function, $u(x, t)$, is the region $[0, L] \times [0, \infty)$. The range, will be the real numbers, $\mathbb{R}$.

How then would we describe the initial heat distribution? We want something that tells us how much heat is at each point, $x$, on our wire at time, $t = 0$. Thus, we have the initial

condition,

$$u(x,0) = f(x).$$

**Question 4.3:** We need some way to describe how the heat in the wire will change as time changes. How do we model heat diffusion?

**Answer:** The heat equation comes from the intuition that heat diffuses in all directions equally. Using this intuition, we begin to model heat flow. For ease, let's pretend that our wire is a close system, so no heat leaves the wire, it just moves around inside it.

1. If the wire is all have the same amount of heat, then what do we expect?

   The heat at each point diffuses equally in all directions. So each point, $x$, sends half of its heat to the left and half to the right. But, at the same time, each point, $x$, receives half of the heat from the points on the left and half of the heat from the points on the right. If all the points have the same amount of heat, $C$, then the change of heat at each point will be

   $$-C + \frac{1}{2}C + \frac{1}{2}C = 0.$$

2. What must happen, then, at points where the temperature goes up?

   Again, we start from the intuition that heat at each point diffuses equally in all directions. So if the temperature goes up at the point, $x$, then $x$ must be receiving more heat that it is giving. Thus, since sends half of its heat to the left and half to the right, it looses all its heat, $C_x$ But, $x$ also receives half of the heat from the points on the left, $\frac{1}{2}C_L$, and half of the heat from the points on the right, $\frac{1}{2}C_R$. Therefore, we can approximate the change of heat at each point will be

   $$-C_x + \frac{1}{2}C_L + \frac{1}{2}C_R > 0.$$

   This means that $C_x < \frac{C_L + C_R}{2}$. That is, **if the amount of heat at $x$ increases, then the amount of heat at $x$ is below the average of the heat at the points around it.**

3. What must happen, then, at points where the temperature goes down?

   Again, we start from the intuition that heat at each point diffuses equally in all directions. So if the temperature goes down at the point, $x$, then $x$ must be losing more heat that it is receiving. Thus, since sends half of its heat to the left and half to the right, it looses all its heat, $C_x$ But, $x$ also receives half of the heat from the points on the left, $\frac{1}{2}C_L$, and half of the heat from the points on the right, $\frac{1}{2}C_R$. Therefore, we can approximate the change of heat at each point will be

   $$-C + \frac{1}{2}C_L + \frac{1}{2}C_R < 0.$$

This means that $C > \frac{C_L + C_R}{2}$. That is, **if the amount of heat at $x$ decreases, then the amount of heat at $x$ is above the average of the heat at the points around it.**

Note that this matches our intuition that if the temperature at a point, $x$, is below the average temperature of the points around it, we expect the temperature at $x$ to increase. Conversely, if the temperature at $x$ is above the average of the temperatures around it, we expect the temperature at $x$ to decrease.

The mathematical concept which captures this sense of "above or below the local average" is **concavity**. Concavity is determined by the second derivative.

The idea of change in time is captured by the derivative with respect to time.

Thus, we can rephrase our intuition by,

1. It the $u_{xx} = 0$, then we expect $u_t = 0$

2. It the $u_{xx} > 0$, then we expect $u_t > 0$

3. It the $u_{xx} < 0$, then we expect $u_t < 0$

There are many equations we could write which have these properties, but perhaps the simplest is

$$Ku_{xx}(x, t) = u_t(x, t)$$

where $K$ is a positive constant.

This is the heat equation. If we let $\partial_x$ represent the operator which takes a partial derivative with respect to $x$ and let $\partial_t$ represent the operator which takes a partial derivative with respect to $t$, then we can do a little algebra and write the heat equation as

$$(K\partial_x^2 - \partial_t)u(x, t) = 0.$$

Let's clean up what we have so far. We want a function, $u(x, t) : [0, L] \times [0, \infty) \to \mathbb{R}$ which satisfies the following conditions:

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x, t) = 0 & \text{for all } (x, t) \in (0, L) \times (0, \infty) \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \end{cases}$$

It turns out that this is not quite well-defined. I need to specify all of the conditions on the boundary of the region $[0, L] \times [0, \infty)$ in order for this problem to have a unique solution. For ease, let's set the side conditions to 0. **That is, we are trying to solve the following problem:**

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x, t) = 0 & \text{for all } (x, t) \in (0, L) \times (0, \infty) \\ u(0, t) = 0 & t > 0 \\ u(L, t) = 0 & t > 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \end{cases}$$

This is a Boundary Value Problem, because we are looking for a function, $u$, which solves a differential equation on the region, $(0, L) \times (0, \infty)$, and satisfies the boundary conditions given above.

**Question 4.4:** What is the structure of the set of solutions to the heat equation?

**Answer:** We answer this the exact same way that we did for linear systems of ODEs. We rewrite the equation, $Ku_{xx}(x,t) = u_t(x,t)$ as

$$(K\partial_x^2 - \partial_t)u(x,t) = 0.$$

**Since $(K\partial_x^2 - \partial_t)$ is a linear operator**– it is left as an exercise to prove it from the definition– **the set of solutions to the heat equation form a vector space.**

**Question 4.5:** How big of a vector space is it? And, how do we get a basis for that vector space?

**Answer:** The vector space is VERY large. It is infinite dimensional, but is so large it does not admit a countable basis. Therefore, we must use a different strategy.

**Question 4.6:** How do we solve BVP involving the heat equation?

**Answer:** The general strategy is the same as before: **find a large enough family of solutions to the heat equation and then compare them with the boundary conditions.** The process of comparing with the boundary conditions will reduce the family of solutions we need to consider so that we can use basis methods developed during out study of Fourier series.

**Question 4.7:** Ok, how do we find solutions to the heat equation? It is a PDE, and PDEs are very hard.

**Answer:** PDEs are very hard. That is why we find solutions to the heat equation by turning it into a pair of ODEs, which we can then solve easily. This trick is called **Separation of Variables**. It is another example of assuming a solution has a particular form, plugging it in to the differential equation, and seeing what comes out.

## 4.2   Separation of Variables

Assume that the function $u(x,y)$ is a solution to the heat equation, $Ku_{xx}(x,t) = u_t(x,t)$,. Let us also assume that $u(x,t)$ can be written as a product of functions,

$$u(x,t) = X(x)T(t).$$

If $u(x,t) = X(x)T(t)$ is a solution to the heat equation, then we can plug it in and get that

$$KX''(x)T(t) = X(x)T'(t).$$

Doing a little algebra to **separate the variables,** we then get that

$$\frac{X''(x)}{X(x)} = \frac{T'(t)}{KT(t)}$$

for all points $x$ and all times $t$. Since the left-hand side only depends upon $x$ and the right-hand side only depends upon $t$, they must be constant in $\Omega$. That is, suppose that $\frac{X''(x)}{X(x)} = \frac{T'(t)}{KT(t)}$ were not constant on $\Omega = [0, L] \times [0, \infty)$. Then, there must be two different points, $(x_1, t_1), (x_2 t_2) \in \Omega$ for which either

$$\frac{T'(t_1)}{KT(t_1)} \neq \frac{T'(t_2)}{KT(t_2)}.$$

or

$$\frac{X''(x_1)}{X(x_1)} \neq \frac{X''(x_2)}{X(x_2)}$$

But, if $\frac{T'(t_1)}{KT(t_1)} = \frac{X''(x)}{X(x)}$ for any $(x, t_1) \in \Omega$ and $\frac{X''(x)}{X(x)} = \frac{T'(t_2)}{KT(t_2)}$ for any $(x, t_1) \in \Omega$, then since both $(x_0, t_1), (x_0, t_2) \in \Omega$ for any $x \in [0, L]$ it must be that

$$\frac{T'(t_1)}{KT(t_1)} = \frac{X''(x_0)}{X(x_0)} = \frac{T'(t_2)}{KT(t_2)}.$$

Thus, this ratio must be constant on all of $\Omega$. We write,

$$\frac{X''(x)}{X(x)} = \frac{T'(t)}{KT(t)} = -\lambda.$$

Now, we can turn these equations into two linear ODEs.

$$\begin{aligned} X''(x) + \lambda X(x) &= 0 \\ T'(t) + \lambda KT(t) &= 0 \end{aligned}$$

Linear ODEs are easy to solve (use 307 techniques). **This is the power of separation of variables: to turn a hard PDE into easy ODEs.**

Separation of variables is a general strategy, and we will use it again for the wave equation and the Laplace equation. Separation of variables does not work on every PDE, but it does work on some very interesting PDEs which we will see in this class.

Now that Separation of variables has done its job, we can get solutions to the heat equation by solving the ODEs, above. This reduces to three cases: when $\lambda > 0$, when $\lambda < 0$, and when $\lambda = 0$.

- Case 1: $\lambda < 0$.

  Solving the temporal ODE, we get $T(t) = ce^{-\lambda K t}$. Solving the spatial ODE, we get $X(x) = ae^{\sqrt{-\lambda}x} + be^{-\sqrt{-\lambda}x}$. Thus, for any $\lambda < 0$, we have the solution

  $$u_\lambda(x,t) = a_\lambda e^{\sqrt{-\lambda}x}e^{-\lambda K t} + b_\lambda e^{-\sqrt{-\lambda}x}e^{-\lambda K t}$$

- Case 2: $\lambda = 0$.

  Solving the temporal ODE, we get $T(t) = c$. Solving the spatial ODE, we get that $X(x) = ax + b$. Thus, we have the solution,

  $$u_0(x,t) = a_0 x + b_0.$$

- Case 3: $\lambda > 0$.

  Solving the temporal ODE, we get $T(t) = ce^{-\lambda K t}$, as before. But, now, solving the spatial ODE, we have $X(x) = a\cos(\sqrt{\lambda}x) + b\sin(\sqrt{\lambda}x)$. Thus, for any $\lambda > 0$, we have the solution

  $$u_\lambda(x,t) = e^{-\lambda K t}(a_\lambda \cos(\sqrt{\lambda}x) + b_\lambda \sin(\sqrt{\lambda}x)).$$

**The collection $\{u_\lambda\}_{\lambda \in \mathbb{R}}$ is our family of solutions to the heat equation.** We should think of it as like a basis for the set of all solutions, but just very, very large.

It would be very hard to try to find a linear combination of the $\{u_\lambda\}$ such that that combination agreed with general boundary conditions. So, we must make our problem easier.

Note that because of Separation of Variables, we have–in some sense– separated how some solutions, $u_\lambda$, behave in space and how they behave in time. In particular, if $u_\lambda = X(x)T(t)$ and there is a point, $x_0$ for which $X(x_0) = 0$, then $u_\lambda(x_0,t) = 0$ for all $t$.

This leads to the idea that if we could impose very simple boundary conditions, like,

$$\begin{cases} u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \end{cases}$$

then, our solution, $u(x,t)$ could only involve $u_\lambda(x,t) = X(x)T(t)$ for $X(x)$ which solve the much simpler BVP,

$$X''(x) + \lambda X(x) = 0; \qquad X(0) = 0 \text{ and } X(L) = 0$$

Thus, by imposing very simple boundary conditions, we can reduce our problem to finding eigenfunctions of a related 2-point BVP. If we can express our initial heat distribution as a linear combination of these eigenfunctions, $X_\lambda(x)$, then that linear combination of the corresponding $u_\lambda(x,t) = X_\lambda(x)T_\lambda(t)$ will be our solution.

## 4.3   Comparing with the Boundary Conditions

To compare this family of solutions with boundary conditions, we need to have some boundary conditions. **The general idea will be that because of the simple boundary conditions we choose, we can reduce our problem to finding eigenfunctions of a related BVP. These eigenfunctions will form a countable basis for the class of functions we care about.**

**Example 4.3:** Since the structure of the set of solutions to the homogeneous BVP determines the structure of the set of solutions to the corresponding non-homogeneous boundary value problems, we investigate the homogeneous case:

$$\begin{cases} Ku_{xx} = u_t & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = 0 & \text{for } x \in [0, L] \end{cases}$$

As mentioned at the end of last section, the collection of solutions $\{u_\lambda\}_{\lambda \in \mathbb{R}}$ is linearly independent. So, if a $u_{\lambda_0}$ is non-zero on the boundaries, there is no non-trivial linear combination of the other $u_\lambda$ which can cancel it out. Thus, **we reduce to only considering $u_\lambda$ which satisfy the boundary conditions at the endpoints**

$$\begin{cases} u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \end{cases}$$

Because of Separation of Variables, this is equivalent to only considering $u_\lambda(x, t) = X_\lambda(x)T_\lambda(t)$ for which $X_\lambda(x)$ satisfies the following 2-point BVP,

$$X''(x) + \lambda X(x) = 0; \qquad X(0) = 0 \text{ and } X(L) = 0.$$

We check, as always, by plugging it in.

- Case 1: $\lambda < 0$.

    $X_\lambda(x) = a_\lambda e^{\sqrt{-\lambda}x} + b_\lambda e^{-\sqrt{-\lambda}x}$. Plugging into our boundary conditions, we have

    $$a + b = 0$$

    $$a(e^{\sqrt{-\lambda}L} - e^{-\sqrt{-\lambda}L}) = 0$$

    where we have used that $a = -b$ from the first equation to simplify the second.

    Since for all $L > 0$, we have that $e^{\sqrt{-\lambda}L} > e^{-\sqrt{-\lambda}L}$, it must be that $a = 0$. Thus, for all $\lambda < 0$, no non-trivial $u_\lambda$ satisfy the boundary condition at infinity.

- Case 2: $\lambda = 0$.

    $X_0(x, t) = a_0 x + b_0$. Plugging in to the end point conditions, we search for $a_0, b_0$ such that

$$b_0 = 0$$

$$a_0 L + b_0 = 0$$

Thus, no non-trivial solutions exist. $a_0 = 0 = b_0$.

- Case 3: $\lambda > 0$.

$X_\lambda(x,t) = a_\lambda cos(\sqrt{\lambda}x) + b_\lambda sin(\sqrt{\lambda}x)$. Plugging in the boundary conditions, we are searching for non-zero constants $a, b$ such that

$$acos(\sqrt{\lambda}0) + bsin(\sqrt{\lambda}0) = 0$$

$$acos(\sqrt{\lambda}L) + bsin(\sqrt{\lambda}L) = 0.$$

The first equation implies $a = 0$, which reduces the second equation to

$$bsin(\sqrt{\lambda}L) = 0.$$

Since we are hoping that $b \neq 0$, we want to know when $sin(\sqrt{\lambda}L) = 0$. This only happens when $\sqrt{\lambda}L = n\pi$ for some $n \in \mathbb{Z}$. Thus, this only happens when $\lambda = \frac{n^2\pi^2}{L^2}$ for some $n \in \mathbb{N}$.

**Thus, the only $\lambda$ for which non-trivial $X_\lambda$ satisfy the 2-point BVP are when $\lambda = \frac{n^2\pi^2}{L^2}$ for some $n \in \mathbb{N}$**

We can enumerate these functions by $n \in \mathbb{N}$ and relabel them

$$X_n(x) = sin(n\frac{\pi}{L}x).$$

**Question 4.8:** Is $\{X_n\}$ a countable basis for $L^2([0, L], \mathbb{R})$?

**Answer:** Yes. A basis must span and be linearly independent. Since every function in $L^2([0, L], \mathbb{R})$ has a Fourier Sine series, it can be written as a linear combination of the functions $\{X_n\}$. This shows that $\{X_n\}$ spans $L^2([0, L], \mathbb{R})$. It is a little trickier to show that the collection $\{X_n\}$ is linearly independent, because it is not an orthogonal basis. However, for our purposes, this does not matter. We shall omit the proof.

Since $\{X_n\}$ is a basis for $L^2([0, L], \mathbb{R})$, for every function, $f \in L^2([0, L], \mathbb{R})$, there exists coefficients $\{b_n\}_{n=1}^\infty$ such that

$$f = \sum_{n=1}^{\infty} b_n sin(n\frac{\pi}{L}x).$$

If we let

$$u_n(x,t) = X_n(x)T_n(t) = e^{-\frac{n^2\pi^2}{L^2}Kt}sin(\frac{n\pi}{L}x),$$

then consider the function

$$u(x,t) = \sum_{n=1}^{\infty} b_n u_n(x,t) = \sum_{n=1}^{\infty} b_n e^{-\frac{n^2\pi^2}{L^2}Kt} sin(\frac{n\pi}{L}x).$$

**Question 4.9:** Is $u(x,t) = \sum_{n=1}^{\infty} b_n u_n(x,t)$ a solution to the heat equation?

**Answer:** Yes. solutions to the heat equation, $(K\partial_x^2 - \partial_t)u(x,t) = 0$, are functions in the kernel of a linear operator, the set of solutions to the heat equation for a vector space. Thus, the linear combination $u(x,t) = \sum_{n=1}^{\infty} b_n u_n(x,t)$ is also a solution to the heat equation.

**Question 4.10:** Does $u(x,t) = \sum_{n=1}^{\infty} b_n u_n(x,t)$ satisfy the boundary conditions?

$$\begin{cases} u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f & \text{for } x \in [0,L] \end{cases}$$

**Answer:** Let's check! We begin with the condition $u(0,t) = 0$ for all $t \geq 0$

$$\begin{aligned} u(0,t) &= \sum_{n=1}^{\infty} b_n u_n(0,t) \\ &= \sum_{n=1}^{\infty} b_n X_n(0) T_n(t) \\ &= \sum_{n=1}^{\infty} b_n 0 \\ &= 0 \end{aligned}$$

Similarly, for the condition that $u(L,t) = 0$ for all $t \geq 0$, we see that

$$\begin{aligned} u(L,t) &= \sum_{n=1}^{\infty} b_n u_n(L,t) \\ &= \sum_{n=1}^{\infty} b_n X_n(L) T_n(t) \\ &= \sum_{n=1}^{\infty} b_n 0 \\ &= 0 \end{aligned}$$

To check the last condition, $u(x,0) = f$ for $x \in [0,L]$, we plug in as always.

$$\begin{aligned} u(x,0) &= \sum_{n=1}^{\infty} b_n u_n(x,0) \\ &= \sum_{n=1}^{\infty} b_n e^{-\frac{n^2\pi^2}{L^2}K0} sin(\frac{n\pi}{L}x) \\ &= \sum_{n=1}^{\infty} b_n sin(\frac{n\pi}{L}x) \\ &= f(x) \end{aligned}$$

because $b_n$ were chosen to be the Fourier Sine series coefficients of $f$.

81

Therefore, if $\{b_n\}$ are the Fourier Sine series coefficients to $\mathcal{F}_{[-L,L]}(f^{odd})$, then

$$u(x,t) = \sum_{n=1}^{\infty} b_n e^{-\frac{n^2\pi^2}{L^2}Kt} sin(\frac{n\pi}{L}x)$$

is the solution to the BVP

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,L] \end{cases}$$

In this example, we are considering the homogeneous case, $u(x,0) = 0$. All the Fourier Sine series coefficients of 0 are 0. This means that **there are no non-trivial solutions to the homogeneous BVP with the heat equation. By our earlier statement, this means non-homogeneous BVP will have unique solutions.**

We now consider two simple non-homogeneous BVP for the heat equation.

**Example 4.4:** Find the solution to the following BVP.

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,L] \end{cases}$$

By exactly the same arguments as above, the solution to the problem is given by the formula

$$u(x,t) = \sum_{n=1}^{\infty} e^{-\frac{n^2\pi^2}{L^2}Kt} b_n sin(\frac{n\pi}{L}x))$$

where the coefficients are given by the formulae

$$b_n = \frac{2}{L} \int_0^L f(x) sin(\frac{n\pi}{L}x) dx.$$

We now solve a slight variation of this type of problem.

**Example 4.5:** Find the solution to following BVP.

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = T_1 & \text{for all } t \geq 0 \\ u(L, t) = T_2 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \end{cases}$$

Now, the endpoints are not zero. We could go back to our collection $\{u_\lambda\}_{\lambda \in \mathbb{R}}$ and compare each of these with our new boundary conditions at the endpoints and infinity. But, there is a smarter way to handle this problem. We instead split the given BVP into two related BVP which we know how to solve.

Consider the two BVPs

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = T_1 & \text{for all } t \geq 0 \\ u(L, t) = T_2 & \text{for all } t \geq 0 \\ u(x, 0) = T_1 - \frac{T_2 - T_1}{L}x & \text{for } x \in [0, L] \end{cases}$$

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) - T_1 + \frac{T_2 - T_1}{L}x & \text{for } x \in [0, L] \end{cases}$$

If we denote the solution to the first by $u_S(x, t)$ and the solution to the second by $u_T(x, t)$, then I claim that $u_S + u_T$ solves the original BVP. Let's check.

$$(K\partial_x^2 - \partial_t)(u_S(x, t) + u_T(x, t)) = 0$$

in the region $\Omega = [0, L] \times [0, \infty)$ because solutions to the heat equation form a vector space. Checking the boundary conditions, we see that

$$\begin{aligned} u_S(0, t) + u_T(0, t) &= T_1 + 0 = T_1 \\ u_S(L, t) + u_T(L, t) &= T_2 + 0 = T_2 \\ u_S(x, 0) + u_T(x, 0) &= T_1 - \frac{T_2 - T_1}{L}x + f(x) - T_1 + \frac{T_2 - T_1}{L}x = f(x) \end{aligned}$$

The important thing about the condition at infinity is that our solutions remain bounded. Since the sum of two bounded functions is bounded, $u_S + u_T$ satisfies that condition, as well. Thus, $u_S(x, t) + u_T(x, t)$ is the unique solution to the given BVP. Now we just need to solve for $u_S$ and $u_T$.

By construction, we already know how to solve for $u_T$ by Example 1.

$$u_T(x, t) = \sum_{n=1}^{\infty} e^{-\frac{n^2 \pi^2}{L^2} Kt} b_n \sin(\frac{n\pi}{L}x))$$

where the coefficients are given by the formulae

$$b_n = \frac{2}{L} \int_0^L (f(x) - T_1 + \frac{T_2 - T_1}{L}x) \sin(\frac{n\pi}{L}x)dx.$$

To solve for $u_S$, we observe that the the the initial heat distribution, $T_1 - \frac{T_2 - T_1}{L}x$ itself solves the equation. Thus,

$$u_S(x, t) = T_1 - \frac{T_2 - T_1}{L}x$$

**Definition 4.6.** We call $u_S$ the **steady-state solution** because it does not depend upon $t$. We call $u_T$ the **transient solution**, because as $t \to \infty$ the function $u_T(x, t)$ decays to zero.

### Reflection Questions

1. What is the formula for the solution to the following BVP for the heat equation?

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x, t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \end{cases}$$

2. Verify that the formula you gave satisfies the above BVP by plugging it into the differential equation and checking the boundary conditions.

3. How did we arrive at the formula in 1., above? Describe the steps in your own words.

4. What is the formula for the solution to the following BVP for the heat equation?

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x, t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = T_1 & \text{for all } t \geq 0 \\ u(L, t) = T_2 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \end{cases}$$

Which parts of the formula give the transient solution? Which parts give the steady-state solution?

5. How did we arrive at the formulae in 3., above? What are the steps that we used?

6. Why are solutions to BVP for the heat equation unique?

7. What does it mean to be an eigenvalue of a linear operator?

8. What does it mean to be the eigenvalue of a BVP?

9. Do you think that $\{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$ is a countable basis for $L^2([a, b], \mathbb{R})$?

## 4.4   Other BVP for the heat equation

In this section, we consider different boundary value problems for the heat equation. In the previous section, we considered heat conduction in a rod with the end-points held at fixed temperatures. Now, we consider what happens if you insulate the end-points. To model insulation, we assume that there is no heat flow our the end-points. Translating this into mathematics, we assume that $\partial_x u(0,t) = 0 = \partial_x u(L,t)$. Thus, the BVP that we shall investigate is,

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ \partial_x u(0,t) = 0 & \text{for all } t \geq 0 \\ \partial_x u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,L] \end{cases}$$

**Because we have changed the type of boundary conditions, we must go back and re-compare with our large family of solutions, $\{u_\lambda\}$.**

However, we can still use the same strategy as before: because our boundary conditions are constant on the end points, the only $u_\lambda$ which will survive the comparison process will be $u_\lambda = X_\lambda(x)T_\lambda(t)$ for which $X_\lambda(x)$ satisfies

$$X''(x) + \lambda X = 0; \qquad X'(0) = 0, X'(L) = 0.$$

That is, we can reduce to finding eigenfunctions of a two-point BVP. To find these eigenfunctions, we find the general solution to the differential equation and then compare with the boundary conditions.

- Case 1: $\lambda < 0$.

  $X_\lambda(x) = a_\lambda e^{\sqrt{-\lambda}x} + b_\lambda e^{-\sqrt{-\lambda}x}$. Plugging into our boundary conditions, we have

  $$a(\sqrt{-\lambda}) - b(\sqrt{-\lambda}) = 0$$

  $$a((\sqrt{-\lambda})e^{\sqrt{-\lambda}L} - (\sqrt{-\lambda})e^{-\sqrt{-\lambda}L}) = 0$$

  where we have used that $a = b$ from the first equation to simplify the second.

  Since for all $L > 0$, we have that $e^{\sqrt{-\lambda}L} > e^{-\sqrt{-\lambda}L})$, it must be that $a = 0$. Thus, $\lambda < 0$, is not an eigenvalue of the BVP because no non-trivial $X_\lambda$ satisfy the BVP.

- Case 2: $\lambda = 0$.

  $X_0(x,t) = a_0 x + b_0$. Plugging in to the end point conditions, we search for $a_0, b_0$ such that

  $$a_0 = 0$$

  $$a_0 = 0$$

  Thus, there are non-trivial solutions. $X_0(x) = b_0$ for any $b_0$ is an eigenfunction of the BVP associated to the eigenvalue, 0.

- Case 3: $\lambda > 0$.

  $X_\lambda(x,t) = a_\lambda cos(\sqrt{\lambda}x) + b_\lambda sin(\sqrt{\lambda}x)$. Plugging in the boundary conditions, we are searching for non-zero constants $a, b$ such that

  $$-a(\sqrt{\lambda})sin(\sqrt{\lambda}0) + b(\sqrt{\lambda})cos(\sqrt{\lambda}0) = 0$$
  $$-a(\sqrt{\lambda})sin(\sqrt{\lambda}L) + b(\sqrt{\lambda})cos(\sqrt{\lambda}L) = 0.$$

  The first equation implies $b = 0$, which reduces the second equation to

  $$-a(\sqrt{\lambda})sin(\sqrt{\lambda}L) = 0.$$

  Since we are hoping to find non-trivial solutions, let's assume that $a \neq 0$. Therefore, we are looking for when $sin(\sqrt{\lambda}L) = 0$. Just as before, this only happens when $\sqrt{\lambda}L = n\pi$ for some $n \in \mathbb{Z}$. Thus, $\lambda = \frac{n^2\pi^2}{L^2}$ for some $n \in \mathbb{N}$.

**Thus the only $\lambda$ for which non-trivial $X_\lambda$ satisfy the 2-point BVP are when $\lambda = \frac{n^2\pi^2}{L^2}$ for some $n \in \mathbb{N}$**

We can enumerate these functions by $n \in \mathbb{N}$ and relabel them

$$X_n(x) = cos(n\frac{\pi}{L}x).$$

**Question 4.11:** Is $\{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$ a countable basis for $L^2([0, L], \mathbb{R})$?

**Answer:** Yes. A basis must span and be linearly independent. Since every function in $L^2([0, L], \mathbb{R})$ has a Fourier Cosine series, it can be written as a linear combination of the functions $\{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$. This shows that $\{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$ spans $L^2([0, L], \mathbb{R})$. It is a little trickier to show that the collection $\{cos(n\frac{\pi}{L}x)\}_{n=0}^{\infty}$ is linearly independent, because it is not an orthogonal basis. We shall omit the proof, but it is a countable basis.

Since $\{cos(n\frac{\pi}{L}x)\}_n$ is a countable basis for $L^2([0, L], \mathbb{R})$, for every function, $f \in L^2([0, L], \mathbb{R})$, there exists coefficients such that

$$f = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n cos(n\frac{\pi}{L}x).$$

This is just $\mathcal{F}_{[-L,L]}(f^{even})$. If we let

$$u_n(x,t) = X_n(x)T_n(t) = e^{-\frac{n^2\pi^2}{L^2}Kt}cos(\frac{n\pi}{L}x),$$

then consider the function

$$u(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n u_n(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n e^{-\frac{n^2\pi^2}{L^2}Kt}cos(\frac{n\pi}{L}x).$$

**Question 4.12:** Is $u(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n u_n(x,t)$ a solution to the heat equation?

**Answer:** Yes. Because a solution to the heat equation, $(K\partial_x^2 - \partial_t)u(x,t) = 0$, is defined by being in the kernel of a linear operator, the set of solutions to the heat equation forms a vector space. Thus, the linear combination $u(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n u_n(x,t)$ is also a solution to the heat equation.

**Question 4.13:** Does $u(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n u_n(x,t)$ satisfy the boundary conditions?

$$\begin{cases} \partial_x u(0,t) = 0 & \text{for all } t \geq 0 \\ \partial_x u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f & \text{for } x \in [0,L] \end{cases}$$

**Answer:** Let's check! We begin with the condition $u(0,t) = 0$ for all $t \geq 0$

$$\begin{aligned} \partial_x u(0,t) &= \sum_{n=1}^{\infty} a_n \partial_x u_n(0,t) \\ &= \sum_{n=1}^{\infty} a_n X_n'(0) T_n(t) \\ &= \sum_{n=1}^{\infty} a_n 0 \\ &= 0 \end{aligned}$$

Similarly, for the condition that $u(L,t) = 0$ for all $t \geq 0$, we see that

$$\begin{aligned} \partial_x u(L,t) &= \sum_{n=1}^{\infty} a_n \partial_x u_n(L,t) \\ &= \sum_{n=1}^{\infty} a_n X_n'(L) T_n(t) \\ &= \sum_{n=1}^{\infty} a_n 0 \\ &= 0 \end{aligned}$$

To check the last condition, $u(x,0) = f$ for $x \in [0,L]$, we plug in as always.

$$\begin{aligned} u(x,0) &= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n u_n(x,0) \\ &= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n e^{-\frac{n^2\pi^2}{L^2}K0} \cos(\frac{n\pi}{L}x) \\ &= \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n \cos(\frac{n\pi}{L}x) \\ &= f(x) \end{aligned}$$

because $a_n$ were chosen to be the Fourier Cosine series coefficients of $f$.

That is, if $\{a_n\}$ are the Fourier Cosine series coefficients to $\mathcal{F}_{[-L,L]}(f^{even})$, then

$$u(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n e^{-\frac{n^2\pi^2}{L^2}Kt} \cos(\frac{n\pi}{L}x)$$

is the solution to the BVP

$$
\begin{cases}
(K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\
\partial_x u(0, t) = 0 & \text{for all } t \geq 0 \\
\partial_x u(L, t) = 0 & \text{for all } t \geq 0 \\
u(x, 0) = f(x) & \text{for } x \in [0, L]
\end{cases}
$$

If we try to solve a BVP for the heat equation with boundary conditions of a more complicated type, it gets very complicated. Sometimes it is possible, but sometimes it is too difficult for the methods we using. We not consider any further BVP for the heat equation.

However, I will mention some very interesting properties of solutions to the heat equation.

1. Solutions to the heat equation are smooth. They have derivatives of all orders.

   This should be somewhat surprizing, since the initial heat distributions, $f$, we considered could be discontinuous. But, no matter how discontinuous our boundary conditions are, inside our region, the solutions are perfectly smooth.

   To see this for the BVP for the heat equation that we considered, we need to recall from the Fourier Series Conceptual Questions that there is a relationship between the smoothness of a function, $f$, and how fast the Fourier Series coefficients of $f$ decay to zero. The more derivatives $f$ has the faster the coefficients decay, and conversely, the faster the Fourier coefficients decay, the more derivatives $f$ must have. In our example,

$$
u(x,t) = \sum_{n=1}^{\infty} b_n u_n(x,t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} a_n e^{-\frac{n^2\pi^2}{L^2}Kt} b_n \cos\left(\frac{n\pi}{L}x\right).
$$

   Thus, for any positive $t > 0$, the collection of "coefficents" looks like

$$
a_n e^{-\frac{n^2\pi^2}{L^2}Kt}.
$$

   Since $n \to \infty$, these decay very, very fast. Thus, $u(x,t)$ is a smooth function in $x$ for every $t > 0$.

2. In this model of heat diffusion, heat travels infinitely fast.

   Obviously, this is kind of a problem. Nothing is supposed to move faster than the speed of light. But in this model, heat moves at infinite speed. What does this mean? It means that if we let our initial heat distribution be, say,

$$
f(x) =
\begin{cases}
0 & x \in [0, 1/3) \\
1 & x \in [1/3, 2/3] \\
0 & x \in (2/3, 1]
\end{cases}
$$

then the solution to

$$\begin{cases} (K\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = (0,1) \times (0,\infty) \\ \partial_x u(0,t) = 0 & \text{for all } t \geq 0 \\ \partial_x u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,1] \end{cases}$$

will be strictly positive for all $(x,t) \in \Omega$. This means that for any $t > 0$, some small amount of heat would have traveled from the middle third to the ends. Thus, it must be infinitely fast.

3. Solutions to BVP for the heat equation satisfy the Maximum Principle.

   This means that
   $$max_\Omega |u(x,t)| = max_{\partial\Omega} |u(x,t)|$$

   This makes a lot of intuitive sense. We don't want heat concentrating. This tells us that the heat diffuses like we expect it to.

4. We can think of the heat equation as an evolution of a graph. That is, for each time, $t > 0$, we can look at the time-slice $u(x,t) = u_t(x)$ as a function of $x$. Thus, we have a family of functions, $\{u_t(x)\}_{t>0}$, As we let $t \to \infty$, we can think of what happens to the graphs of these functions. It turns out that this describes a way to reduce the arc length of the graph in a "fastest" way.

## 4.5   The wave equation

In this section, we will consider BVP for the wave equation. That is, we imagine that we have a string, which we give some initial position and some initial velocity – by, say, plucking it– and we want to know how the string will vibrate as time goes on.

**Question 4.14:** How do we model this with mathematics?

**Answer:** We let $x$ be the spatial variable. Suppose that our string is of length $L > 0$. Let $t$ represent time. Our region, then, is

$$\Omega = [0,L] \times [0,\infty)$$

and we want a function, $u(x,t) : [0,L] \times [0,\infty) \to \mathbb{R}$, which tells us the position of the point $x$ on the string at time $t$.

Our boundary conditions will be some conditions of the endpoints, say,

$$\begin{cases} u(0,t) & = g_0(t) \\ u(L,t) & = g_L(t) \end{cases}$$

and functions which describe the initial displacement and initial velocity of the string,

$$\begin{cases} u(x,0) & = f(x) \\ \partial_t u(x,0) & = g(x) \end{cases}$$

The differential equation we will have will be the wave equation.

**Question 4.15:** What is the wave equation?

**Answer:** The wave equation is a Partial Differential Equation (PDE) which tries to capture the behavior of waves. We expect waves to oscillate. That is, a particular point on a string will fluctuate back and forth. If we think of concavity, it is not necessary that a point at which the string is concave down will move down, but the forces will pull it down. This means that the acceleration will be the same sign as the concavity. We assume that the relationship is as simple as possible and express this in mathematics by the PDE,

$$\alpha^2 \partial_x^2 u(x,t) = \partial_t^2 u(x,t).$$

**Question 4.16:** What is the structure of the set of solutions to the wave equation?

**Answer:** As you might expect, we can re-write the wave equation as

$$(\alpha^2 \partial_x^2 - \partial_t^2) u(x,t) = 0.$$

Thus, solutions are functions in the kernel of a linear operator. As we have seen, this means that solutions form a vector space.

**Question 4.17:** How big is this vector space? Can we find a basis?

**Answer:** Just as for the heat equation, this vector space is very, very large. It has a basis, but not a countable one. We will need to us the same tricks we used for the heat equation to solve BVP for the wave equation.

**Question 4.18:** How do we solve the wave equation?

**Answer:** Just as in the heat equation, we assume very simple (constant) boundary conditions that allow us to reduce a BVP for the wave equation to finding eigenfunctions for a related two-point BVP.

How did we do this? First, we used Separation of Variables to transform the linear PDE into solvable linear ODEs. Solving these ODEs, we can generate a family of solutions, $\{u_\lambda\}$. By the linear independence of these solutions, we can reduce to comparing these solutions to end-point conditions (i.e., reduce to solving for eigenfunctions of a two-point BVP). Then, we use Fourier series to sum up the solutions to match the other boundary conditions. We make this explicit, below.

Separation of Variables goes as follows, assume that $u(x,t) = X(x)T(t)$ is a solution to the wave equation. Then,

$$\alpha^2 X''(x)T(t) = X(x)T''(t).$$

Doing a little algebra to **separate the variables,** we then get that

$$\frac{X''(x)}{X(x)} = \frac{T''(t)}{\alpha^2 T(t)}$$

for all points $x$ and all times $t$. Just as in the heat equation, since the left-hand side only depends upon $x$ and the right-hand side only depends upon $t$, they must be constant. Thus, we write

$$\frac{X''(x)}{X(x)} = \frac{T''(t)}{\alpha^2 T(t)} = -\lambda$$

Now we can turn these equations into two linear ODEs.

$$
\begin{aligned}
X''(x) + \lambda X(x) &= 0 \\
T''(t) + \lambda \alpha^2 T(t) &= 0
\end{aligned}
$$

Now, we solve these linear ODEs.

- Case 1: $\lambda < 0$.

  Solving the temporal ODE, we get $T(t) = ae^{\alpha\sqrt{-\lambda}t} + be^{-\alpha\sqrt{-\lambda}t}$. Solving the spatial ODE, we get $X(x) = ce^{\sqrt{-\lambda}x} + de^{-\sqrt{-\lambda}x}$. Thus, for any $\lambda < 0$, we have the solution

  $$u_\lambda(x,t) = (c_\lambda e^{\sqrt{-\lambda}x} e^{-\lambda Kt} + d_\lambda e^{-\sqrt{-\lambda}x})(a_\lambda e^{\alpha\sqrt{-\lambda}t} + b_\lambda e^{-\alpha\sqrt{-\lambda}t})$$

- Case 2: $\lambda = 0$.

  Solving the temporal ODE, we get $T(t) = at + b$. Solving the spatial ODE, we get that $X(x) = cx + d$. Thus, we have the solution,

  $$u_0(x,t) = (a_0 t + b_0)(c_0 x + d_0).$$

- Case 3: $\lambda > 0$.

  Solving the temporal ODE, we get $T(t) = a\cos(\alpha\sqrt{\lambda}t) + b\sin(\alpha\sqrt{\lambda}t)$, as before. But, now, solving the spatial ODE, we have $X(x) = c\cos(\sqrt{\lambda}x) + d\sin(\sqrt{\lambda}x)$. Thus, for any $\lambda > 0$, we have the solution

  $$u_\lambda(x,t) = (a_\lambda \cos(\alpha\sqrt{\lambda}t) + b_\lambda \sin(\alpha\sqrt{\lambda}t))(c_\lambda \cos(\sqrt{\lambda}x) + d_\lambda \sin(\sqrt{\lambda}x)).$$

These three formulas give us a very large class of solutions to the wave equation.

**Example 4.6.** We consider the homogeneous case,

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = 0 & \text{for } x \in [0,L] \\ \partial_t u(x,0) = 0 & \text{for } x \in [0,L] \end{cases}$$

By the linear independence of the solutions, $u_\lambda = X_\lambda T_\lambda$, above, we can determine that if the solution to the BVP for the wave equation is a linear combination of the $u_\lambda$ then the $u_\lambda$ which have non-zero coefficients must satisfy

$$\begin{cases} u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \end{cases}$$

That is, $u$ can only be a linear combination of $u_\lambda = X_\lambda T_\lambda$ for $X_\lambda$ which satisfy

$$X''(x) + \lambda X(x) = 0; \qquad X(0) = 0, X(L) = 0$$

But, this is the same problem we solved for the heat equation. Thus, we know that the only eigenfunctions of the two-point BVP are

$$X_n(x) = sin(n\frac{\pi}{L}x).$$

Thus, we have the family $u_n(x,t) = sin(n\frac{\pi}{L}x)(a_n cos(\alpha \frac{n\pi}{L}t) + b_n sin(\alpha \frac{n\pi}{L}t))$. We can assume that

$u(x,t) = \sum_{n=1}^{\infty} sin(n\frac{\pi}{L}x)(a_n cos(\alpha \frac{n\pi}{L}t) + b_n sin(\alpha \frac{n\pi}{L}t))$

Now we need to deal with our initial displacement and initial velocity conditions. We plug in to check them.

$$\begin{aligned} 0 &= u(x,0) \\ &= \sum_{n=1}^{\infty} sin(n\frac{\pi}{L}x)(a_n cos(\alpha \frac{n\pi}{L}0) + b_n sin(\alpha \frac{n\pi}{L}0)) \\ &= \sum_{n=1}^{\infty} a_n sin(n\frac{\pi}{L}x) \end{aligned}$$

By Fourier series results, we know that the coefficients in the Fourier Sine series of the zero function are all zero. Thus, $a_n = 0$ for all $n$. Similarly,

$$\begin{aligned} 0 &= \partial_t u(x,0) \\ &= \sum_{n=1}^{\infty} sin(n\frac{\pi}{L}x)(-a_n \alpha \frac{n\pi}{L} sin(\alpha \frac{n\pi}{L}0) + b_n \alpha \frac{n\pi}{L} cos(\alpha \frac{n\pi}{L}0)) \\ &= \sum_{n=1}^{\infty} b_n \alpha \frac{n\pi}{L} sin(n\frac{\pi}{L}x) \end{aligned}$$

Again, this implies that $b_n = 0$ for all $n$. Thus, the only solution to the homogeneous equation is the trivial solution.

**Example 4.7.** Find the solution to the following BVP,

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \\ \partial_t u(x, 0) = g(x) & \text{for } x \in [0, L] \end{cases}$$

Since the end-point conditions are the same, by identical argument as before, we have that

$$u(x,t) = \sum_{n=1}^{\infty} \sin(n\frac{\pi}{L}x)(a_n \cos(\alpha\frac{n\pi}{L}t) + b_n \sin(\alpha\frac{n\pi}{L}t)).$$

To check against our initial displacement and initial velocity conditions, we simply plug in.

$$\begin{aligned} f(x) &= u(x, 0) \\ &= \sum_{n=1}^{\infty} \sin(n\frac{\pi}{L}x)(a_n \cos(\alpha\frac{n\pi}{L}0) + b_n \sin(\alpha\frac{n\pi}{L}0)) \\ &= \sum_{n=1}^{\infty} a_n \sin(n\frac{\pi}{L}x) \end{aligned}$$

Thus, $\{a_n\}$ are the coefficients to $\mathcal{F}_{[-L,L]}(f^{odd})$. That is, they are the Fourier Sine series coefficients of $f$.

Checking our other condtition,

$$\begin{aligned} g(x) &= \partial_t u(x, 0) \\ &= \sum_{n=1}^{\infty} \sin(n\frac{\pi}{L}x)(-a_n \alpha\frac{n\pi}{L} \sin(\alpha\frac{n\pi}{L}0) + b_n \alpha\frac{n\pi}{L} \cos(\alpha\frac{n\pi}{L}0)) \\ &= \sum_{n=1}^{\infty} b_n \alpha\frac{n\pi}{L} \sin(n\frac{\pi}{L}x) \end{aligned}$$

Thus, $\{b_n \alpha\frac{n\pi}{L}\}$ are the coefficients to $\mathcal{F}_{[-L,L]}(g^{odd})$. That is, they are the Fourier Sine series coefficients of $g$.

Using our Fourier series formulae, then, we have that,

$$a_n = \frac{2}{L} \int_0^L f(x) \sin(n\frac{\pi}{L}x) dx$$

and,

$$b_n = \frac{2}{\alpha n\pi} \int_0^L g(x) \sin(n\frac{\pi}{L}x) dx.$$

**Example 4.8.** Solve the following BVP.

$$\begin{cases} (4\partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0, \pi] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(\pi, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, \pi] \\ \partial_t u(x, 0) = g(x) & \text{for } x \in [0, \pi] \end{cases}$$

93

where $f(x) = \begin{cases} x & [0, \pi/2] \\ 0 & (\pi/2, \pi] \end{cases}$ and $g(x) = 2sin(2x) - 10sin(20x)$.

By the previous example, and our choice of $\alpha = 2$ and $L = \pi$, we already know that the solution, $u(x, t)$ is given by the formulae,

$$u(x, t) = \sum_{n=1}^{\infty} a_n sin(nx)cos(2nt) + \sum_{n=1}^{\infty} b_n sin(nx)sin(2nt).$$

where the coefficients $\{a_n\}$ and $\{b_n\}$ satisfy the equations,

$$a_n = \frac{2}{\pi} \int_0^{\pi} f(x)sin(nx)dx$$

and

$$b_n = \frac{2}{2n\pi} \int_0^{\pi} g(x)sin(nx)dx.$$

Now we need only calculate the coefficients. We calculate the $b_n$, first. Since $g(x)$ is it's own Fourier Sine series, we can read off the coefficients,

$$b_2 = 2, b_{20} = -10$$

and $b_n = 0$ for all other $n$.

To calculate $a_n$, we must compute the integrals.

$$\begin{aligned} a_n &= \frac{2}{\pi} \int_0^{\pi} f(x)sin(nx)dx \\ &= \frac{2}{\pi} \int_0^{\pi/2} xsin(nx)dx \\ &= \frac{-2}{n\pi}cos(nx)x|_0^{\pi/2} - \int_0^{\pi/2} \frac{-2}{n\pi}cos(nx)dx \\ &= \frac{2}{n^2\pi}sin(nx)|_0^{\pi/2} \\ &= \frac{2}{n^2\pi}sin(\frac{n\pi}{2}) \end{aligned}$$

Thus, the solution is

$$u(x, t) = 2sin(2x)sin(4t) - 10sin(20x)sin(40t) + \sum_{n=1}^{\infty} \frac{2}{n^2\pi}sin(\frac{n\pi}{2})sin(nx)cos(2nt).$$

**Question 4.19:** What if we want to know the position of a certain point at a particular time? **How do we compute** $u(x, t)$? These infinite series seem very ugly.

**Answer:** There is a different way to compute solutions to the wave equation, one which is much easier to compute than these Fourier series methods we have been using. This other way of getting solutions is called the **method of characteristics.**

## 4.6 The method of characteristics

To get solutions to BVP for the wave equation using the method of characteristics, we will use a familiar strategy. The method of characteristics itself will give us a large class of solutions to the wave equation. Then, we use the boundary conditions to compare with the solutions to determine the particular solution to the BVP.

As for the method of characteristics itself, it is another "guess and check" method. That is, **we guess the form of a solution and then plug it in to check what conditions functions in that form must solve in order to be a solution.** In this particular instance, we assume that our solution, $u(x, t)$, can be written as

$$u(x, t) = h(\gamma(x, t))$$

for some $\gamma : [0, L] \times [0, \infty) \to \mathbb{R}$ and $h : \mathbb{R} \to \mathbb{R}$. Let's consider what this means. We are assuming that there is a function, $\gamma$, such that the value of $u(x, t)$ only depends upon the level sets of $\gamma$. That is, if $(x_1, t_1)$ and $(x_2, t_2)$ are two points at two times such that

$$\gamma(x_1, t_1) = \gamma(x_2, t_2)$$

then

$$
\begin{aligned}
u(x_1, t_1) &= h(\gamma(x_1, t_1)) \\
&= h(\gamma(x_2, t_2)) \\
&= u(x_2, t_2)
\end{aligned}
$$

Thus, the position or height of a point at a given time only depends upon the level sets of $\gamma$.

Now, we plug in this guess and see what happens. If $u(x, t) = h(\gamma(x, t))$ is a solution to the wave equation, then $(\alpha^2 \partial_x^2 - \partial_t^2)u = 0$. Expanding this, we get that,

$$\alpha^2(h''(\gamma(x, t))(\partial_x\gamma(x, t))^2 + h'(\gamma(x, t))\partial_x^2\gamma(x, t)) - h''(\gamma(x, t))(\partial_t\gamma(x, t))^2 + h'(\gamma(x, t))\partial_t^2\gamma(x, t)) = 0$$

Unlike separation of variables, we have not traded our PDE for a few ODE. Here, we have taken a PDE and seemingly made it more complicated. However, here is where we make an assumption about $\gamma$ which will simplify things greatly.

What happens if $\gamma(x, t) = a + bx + ct$? If $\gamma(x, t) = a + bx + ct$, then all of it's second order partial derivatives are zero. Thus, our equation simplifies to,

$$\alpha^2 h''(\gamma(x, t))(\partial_x\gamma(x, t))^2 - h''(\gamma(x, t))(\partial_t\gamma(x, t))^2 = 0$$

This is much easier to deal with. Since both terms have $h''$ as a factor, we can rewrite this as

$$\alpha^2(\partial_x\gamma(x, t))^2 = (\partial_t\gamma(x, t))^2$$

This is still a PDE, but since we have assumed that $\gamma = a + bx + ct$, we can simply plug in and solve for the coefficients, $a, b, c$.

Since $\partial_x\gamma(x, t) = b$ and $\partial_t\gamma(x, t) = c$, we simply have the equation,

$$\alpha^2 b^2 = c^2$$

Since $a$ can be anything, we choose $a = 0$ . Choosing $b = 1$ gives that $\gamma(x, t) = x \pm \alpha t$. That is, both $\gamma(x, t) = x + \alpha t$ and $\gamma(x, t) = x - \alpha t$ satisfy the equation.

Thus, our large class of solutions are $h(x + \alpha t)$ and $h(x - \alpha t)$ for any twice-differentiable function $h$. This is a very large class of solutions.

**Question 4.20:** How do the functions $h(x + \alpha t)$ and $h(x - \alpha t)$ behave?

**Answer:** Of course, this depends upon the function $h$, but lets consider what happens when we fix $x$. These functions describe a change of coordinates by translation. That is,

$$f(y) = h(x + \alpha t)$$

if $y = x + \alpha t$. So, for any fixed $t$, the graph of $h(x + \alpha t)$ is the graph of $h(x)$, but translated to the left by $\alpha t$.

Similarly, the graph of $h(x + \alpha t)$ is the graph of $h(x)$, but translated to the right by $\alpha t$.

**Question 4.21:** Ok. How do we get solutions to a BVP for the wave equation, then?

**Answer:** As mentioned before, we need to compare solutions to the wave equation to the boundary conditions. So, let's get ourselves some boundary conditions.

**Example 4.9:** Let's get solutions to the following BVP for the wave equation.

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \\ \partial_t u(x, 0) = 0 & \text{for } x \in [0, L] \end{cases}$$

**To compare with the boundary conditions of a BVP for the wave equation, we assume that the solution is of the form**

$$u(x, t) = h_1(x + \alpha t) + h_2(x - \alpha t).$$

We now plug into the boundary conditions and get that $h_1(x \pm \alpha t)$ must satisfy

$$\begin{cases} h_1(\alpha t) + h_2(-\alpha t) = 0 & \text{for all } t \geq 0 \\ h_1(L + \alpha t) + h_2(L - \alpha t) = 0 & \text{for all } t \geq 0 \\ h_1(x) + h_2(x) = f(x) & \text{for } x \in [0, L] \\ \alpha(h_1'(x) - h_2'(x)) = 0 & \text{for } x \in [0, L] \end{cases}$$

We begin with the last condition. Since $h_1' = h_2'$, integrating each side, we see that

$$h_1(x) = h_2(x) + c$$

for $x \in [0, L]$. For convenience, we choose $c = 0$. By our next condition, though, we have that $h_1(x) + h_2(x) = f(x)$ for $x \in [0, L]$. Thus, substituting in, we see that

$$h_1(x) = h_2(x) = \frac{1}{2}f(x)$$

for $x \in [0, L]$.

Our first condition, then, that $h_1(\alpha t) + h_2(-\alpha t) = 0$ for all $t \geq 0$ tells us that

$$h_2(-\alpha t) = -\frac{1}{2}f(\alpha t).$$

Chasing this through, this means that both $h_1$ and $h_2$ are odd functions.

We consider the last condition, that $h_1(L + \alpha t) + h_2(L - \alpha t) = 0$ for all $t \geq 0$. We do a change of variable, $t = \frac{L+y}{\alpha}$. This gives that,

$$h_1(2L + y) + h_2(-y) = h_1(2L + y) - h_1(y) = 0$$

Since this holds for all $y$, we have that $h_1$ is $2L-$periodic. Similarly, we can then show that $h_2$ must also be $2L-$periodic.

**Thus, $h_1$ and $h_2$ are the odd, $2L-$periodic extensions of $\frac{1}{2}f(x)$. Call this, $\tilde{f}(x)$.**
Thus, our solution, $u(x, t)$, is given by the formula,

$$u(x, t) = \frac{1}{2}[\tilde{f}(x + \alpha t) + \tilde{f}(x - \alpha t)].$$

This solution has a clear physical interpretation. The solution is the average of two copies of the off, $2L-$periodic extension of the initial displacement. One copy which travels left at speed $\alpha t$ and one copy which travels right at speed $\alpha t$. This captures the intuition that the wave should travel in both directions and that the solution should be the sum of these wave forms.

**Example 4.10:** Let's get solutions to the following BVP for the wave equation.

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x, t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = 0 & \text{for all } t \geq 0 \\ u(L, t) = 0 & \text{for all } t \geq 0 \\ u(x, 0) = 0 & \text{for } x \in [0, L] \\ \partial_t u(x, 0) = g(x) & \text{for } x \in [0, L] \end{cases}$$

We proceed by exactly the same process; we assume that $u(x,t) = h_1(x+\alpha t) + h_2(x-\alpha t)$ and then plug this into the boundary conditions. This gives that $h_1$ and $h_2$ must satisfy

$$\begin{cases} h_1(\alpha t) + h_2(-\alpha t) = 0 & \text{for all } t \geq 0 \\ h_1(L + \alpha t) + h_2(L - \alpha t) = 0 & \text{for all } t \geq 0 \\ h_1(x) + h_2(x) = 0 & \text{for } x \in [0, L] \\ \alpha(h_1'(x) - h_2'(x)) = g(x) & \text{for } x \in [0, L] \end{cases}$$

Differentiating the third condition, $h_1(x) + h_2(x) = 0$, gives that $h_1' + h_2' = 0$. Plugging this into the fourth condition, we see that

$$h_1'(x) = \frac{1}{2\alpha} g(x)$$

and

$$h_2'(x) = \frac{-1}{2\alpha} g(x)$$

for $x \in [0, L]$. Integrating, we have that for $x \in [0, L]$,

$$h_1(x) = \frac{1}{2\alpha} \int_0^x g(s)ds$$

and

$$h_2(x) = \frac{1}{2\alpha} \int_x^0 g(s)ds.$$

Again, choosing $t = \frac{L+x}{\alpha}$ in the second conditions give us

$$h_1(2L + x) + h_2(x) = 0, \quad \forall x \geq 0.$$

Since $h_1(x) + h_2(x) = 0$, this implies that both $h_1$ and $h_2$ must be $2L-$periodic.

The first condition, then, tells us that $h_1(x)$ and $h_2(x)$ are odd. These conditions force us to extend $g$ to be odd and $2L-$periodic. Call this odd, $2L-$periodic extension $\tilde{g}$. The formula becomes,

$$u(x,t) = h_1(x + \alpha t) + h_2(x + \alpha t) = \frac{1}{2\alpha} \int_{x-\alpha t}^{x+\alpha t} \tilde{g}(s)ds$$

**Example 4.11:** Solve the following BVP,

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0, L] \times [0, \infty) \\ u(0, t) = l_1 & \text{for all } t \geq 0 \\ u(L, t) = l_2 & \text{for all } t \geq 0 \\ u(x, 0) = f(x) & \text{for } x \in [0, L] \\ \partial_t u(x, 0) = g(x) & \text{for } x \in [0, L] \end{cases}$$

We should view this BVP as really the sum of three different BVPs for the wave equation. That is, if we let $u_L(x)$ be the solution to

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = l_1 & \text{for all } t \geq 0 \\ u(L,t) = l_2 & \text{for all } t \geq 0 \\ u(x,0) = l_1 + \frac{(l_2 - l_1)}{L}x & \text{for } x \in [0,L] \\ \partial_t u(x,0) = 0 & \text{for } x \in [0,L] \end{cases}$$

and $u_f(x,t)$ be the solution to

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) - l_1 + \frac{(l_2 - l_1)}{L}x & \text{for } x \in [0,L] \\ \partial_t u(x,0) = 0 & \text{for } x \in [0,L] \end{cases}$$

and $u_g(x,t)$ be the solution to

$$\begin{cases} (\alpha^2 \partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = 0 & \text{for all } t \geq 0 \\ u(L,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = 0 & \text{for } x \in [0,L] \\ \partial_t u(x,0) = g(x) & \text{for } x \in [0,L] \end{cases}$$

Then, just as for the heat equation, $u(x,t) = u_L(x) + u_f(x,t) + u_g(x,t)$. The previous two examples gave equations for $u_f(x,t)$ and $u_g(x,t)$, respectively. To finish the solution, we need only find $u_L(x)$.

However, just as in the heat equation, $u_L(x) = l_1 + \frac{(l_2 - l_1)}{L}x$ satisfies the BVP. This completes the solution for the wave equation by the method of characteristics.

Together, these formulae give what is called **D'Alembert's solution to BVP for the wave equation.**

**Question 4.22:** Do we need $f$ and $g$ to be twice-differentiable for D'Alembert's solutions?

**Answer:** We started off with the assumption for the method of characteristics that $h$ was twice differentiable. But, interestingly, **the solutions that we found do not depend**

**upon being able to differentiate $f$ or $g$. Thus, we can extend these solutions to non-differentiable functions. In fact, we can even extend these solutions to work for discontinuous solutions!**

Thus, D'Alembert's solutions give us the solutions for even discontinuous initial displacement and velocity functions.

Here are some interesting facts about the wave equation to contrast it from the heat equation.

1. Solutions to the wave equation are not smooth.

   This is easy to see from D'Alembert's solution. If say $u(x,t) = \frac{1}{2}[\tilde{f}(x+\alpha t) + \tilde{f}(x-\alpha t)]$, then we see that if $f$ has a discontinuity, then so will $\tilde{f}$. Therefore, this discontinuity will simply be propagated left and right for eternity. There is no "improvement" of the initial displacement and velocity functions.

   We can also see this from the Fourier Series solutions. Recall from the Fourier Series Conceptual Questions that there is a relationship between the smoothness of a function, $f$, and how fast the Fourier Series coefficients of $f$ decay to zero. The more derivatives $f$ has the faster the coefficients decay, and conversely, the faster the Fourier coefficients decay, the more derivatives $f$ must have. In our example, say,

   $$u(x,t) = \sum_{n=1}^{\infty} a_n \sin(\frac{n\pi}{L}x)\cos(\alpha\frac{n\pi}{L}t).$$

   Thus, for any positive $t > 0$, the collection of "coefficents" looks like

   $$a_n \cos(\alpha\frac{n\pi}{L}t).$$

   So, if $t = \frac{2L}{\alpha}$, then all the coefficient are the same as for $t = 0$. Thus, it is clear that we cannot impose any rate of decay based upon $t$, since the coefficients merely oscillate in $t$.

2. In this model of the wave diffusion, a wave travels travels at a finite speed. It does not travel infinitely fast.

   Again, this is easiest to see from D'Alembert's solutions.

   $$f(x) = \begin{cases} 0 & x \in [0, 1/3) \\ 1 & x \in [1/3, 2/3] \\ 0 & x \in (2/3, 1] \end{cases}$$

   then the solution to

   $$\begin{cases} (\partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,1] \times [0,\infty) \\ \partial_x u(0,t) = 0 & \text{for all } t \geq 0 \\ \partial_x u(1,t) = 0 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,1] \\ \partial_t u(x,0) = 0 & \text{for } x \in [0,1] \end{cases}$$

will be $u(x,t) = \frac{1}{2}[\tilde{f}(x+t) + \tilde{f}(x-t)]$ for $\tilde{f}$ the odd, 2−periodic extension of $f$. The point $x = 1/6$ on the string does not move until time $t = 1/6$. The wave takes that long to get to it.

## Reflection Questions:

1. Are solutions to the non-homogeneous BVP for the wave equation unique? Why?

2. What are the formulae for solving a BVP for the wave equation,

$$\begin{cases} (\alpha^2\partial_x^2 - \partial_t^2)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = l_1 & \text{for all } t \geq 0 \\ u(L,t) = l_2 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,L] \\ \partial_t u(x,0) = g(x) & \text{for } x \in [0,L] \end{cases}$$

   using Fourier Series methods?

3. What are the formulae for solving a BVP for the wave equation,

$$\begin{cases} (\alpha^2\partial_x^2 - \partial_t)u(x,t) = 0 & \text{in } \Omega = [0,L] \times [0,\infty) \\ u(0,t) = l_1 & \text{for all } t \geq 0 \\ u(L,t) = l_2 & \text{for all } t \geq 0 \\ u(x,0) = f(x) & \text{for } x \in [0,L] \\ \partial_t u(x,0) = g(x) & \text{for } x \in [0,L] \end{cases}$$

   using the method of characteristics?

4. What are the steps that we used in the method of characteristics to get solutions to the wave equation? Describe them in your own words.

5. What does it mean to be the eigenvalue of a linear operator? What does it mean to be the eigenfunction of a linear operator? (This was in the beginning of section 4.)