

Geometric Deep Learning

MLSS Skoltech, Moscow, August 2019

Michael Bronstein

Mathematical Background

1 Vector spaces

V is a *vector space* over a field \mathbb{F} (typically, $\mathbb{F} = \mathbb{R}$ or \mathbb{C}) with binary operations $+$: $V \times V \rightarrow V$ (*vector addition*) and \cdot : $\mathbb{F} \times V \rightarrow V$ (*vector-by-scalar multiplication*) if for any $u, v, w \in V$ and $\alpha, \beta \in \mathbb{F}$ we have:

1. $u + (v + w) = (u + v) + w$ *associativity of $+$*
2. $u + v = v + u$ *commutativity of $+$*
3. $\exists! 0 \in V$ s.t. $u + 0 = u$ *identity element of $+$*
4. $\exists! (-u) \in V$ s.t. $u + (-u) = 0$ *inverse element of $+$*
5. $\alpha \cdot (u + v) = \alpha \cdot u + \alpha \cdot v$ *distributivity of \cdot w.r.t. vector addition*
 $(\alpha + \beta) \cdot v = \alpha \cdot v + \beta \cdot v$ *distributivity of \cdot w.r.t. scalar addition*
6. $\alpha \cdot (\beta \cdot v) = (\alpha \cdot \beta) \cdot v$ *compatibility of \cdot with scalar multiplication*
7. $\exists! 1 \in \mathbb{F}$ s.t. $1 \cdot u = u$ *identity element of \cdot*

Remarks. Note that we do not include closure $u + v \in V$ and $\alpha \cdot v \in V$ as a separate axioms since we defined $+$: $V \times V \rightarrow V$ and \cdot : $\mathbb{F} \times V \rightarrow V$.

Remarks. Notation sometimes can be confusing and therefore used with care:

1. The same notation is used for scalar addition $\alpha + \beta$ and vector addition $u + v$. It should be understood from context which addition is meant.
2. The same notation is used for scalar-by-scalar multiplication $\alpha \cdot \beta$ and vector-by-scalar multiplication $\alpha \cdot u$.
3. When no confusion arises, very often the vector by scalar multiplication is denoted as αu for brevity.
4. The zero vector $0 \in V$ (identity element of vector addition) should not be confused with the zero scalar $0 \in \mathbb{R}$ (identity element of scalar addition), very often denoted in the same way.

Examples

1. *Vectors*: $\mathbb{R}^n = \{(v_1, \dots, v_n) : v_i \in \mathbb{R} \ \forall i = 1, \dots, n\}$ with $u + v = (u_1 + v_1, \dots, u_n + v_n)$
2. *Functions*: $\mathcal{F}(\Omega) = \{f : \Omega \rightarrow \mathbb{R}\}$ with $(f + g)(x) = f(x) + g(x)$.

2 Norm

Given a vector space V over \mathbb{C} , a *norm* is a function $\|\cdot\| : V \rightarrow \mathbb{R}$ satisfying for any $u, v \in V$ and $\alpha \in \mathbb{C}$:

1. $\|\alpha u\| = |\alpha| \|u\|$ *positive homogeneity*
2. $\|u + v\| \leq \|u\| + \|v\|$ *triangle inequality*
3. $\|u\| = 0 \Rightarrow u = 0$

$(V, \|\cdot\|)$ is called a *normed (vector) space*. Intuitively, the norm measures the length of a vector.

Remark. The following properties (often listed as part of axiomatic definition of the norm) are in fact consequences of the above definition:

1. $\|0\| = \|0 \cdot u\| \stackrel{(1)}{=} |0| \|u\| = 0$, i.e. property (3) is iff: $\|u\| = 0 \iff u = 0$.
2. $\|u\| \geq 0$.

Remark. A function satisfying only properties (1)+(2) is called a *seminorm* and sometimes denoted by $|u|$. A seminorm can assign zero values to vectors that are not necessarily zero vectors.

Examples

1. *Euclidean norm on \mathbb{C}^n* : $\|u\|_2 = \sqrt{\sum_{i=1}^n |u_i|^2}$
2. *L_p -norm on \mathbb{C}^n* : $\|u\|_p = (\sum_{i=1}^n |u_i|^p)^{1/p}$, in particular
 L_1 -norm: $\|u\|_1 = \sum_{i=1}^n |u_i|$
 L_∞ -norm: $\|u\|_\infty = \max\{|u_1|, \dots, |u_n|\}$
3. *L_p -norm on $\mathcal{F}(\Omega)$* : $\|f\|_p = (\int_\Omega |f(x)|^p dx)^{1/p}$, where dx denotes the appropriately defined volume element on Ω .

3 Inner product

Given a vector space V over \mathbb{C} , an *inner product* is a function $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{C}$ satisfying for any $u, v, w \in V$ and $\alpha \in \mathbb{C}$:

1. $\langle u, v \rangle = \overline{\langle v, u \rangle}$ *conjugate (Hermitian) symmetry*

2. $\langle \alpha u, v \rangle = \alpha \langle u, v \rangle$ *linearity*
 $\langle u + w, v \rangle = \langle u, v \rangle + \langle w, v \rangle$
3. $\langle u, u \rangle \geq 0$
 $\langle 0, 0 \rangle = 0$

$(V, \langle \cdot, \cdot \rangle)$ is called an *inner product space*. Intuitively, the inner product can be related to the angle between two vectors (see below).

Remark. The additional property $\langle u, \alpha v \rangle = \overline{\alpha} \langle u, v \rangle$ follows from (1)+(2).

Examples

1. Real vectors \mathbb{R}^n : $\langle u, v \rangle = \sum_{i=1}^n u_i v_i = v^\top u$
2. Complex vectors \mathbb{C}^n : $\langle u, v \rangle = \sum_{i=1}^n u_i \overline{v_i} = v^* u$
3. Real matrices: $\langle A, B \rangle = \text{trace}(AB^\top)$
4. Square-integrable functions $L^2(\Omega)$: $\langle f, g \rangle = \int_{\Omega} f(x) \overline{g(x)} dx$
5. Square-summable real sequences ℓ^2 : $\langle x, y \rangle = \sum_{i \geq 1} x_i y_i$

Relation to norms

1. Inner product naturally defines a norm $\|u\| = \langle u, u \rangle^{1/2}$. Such a norm satisfies the *Cauchy-Schwarz-(Bunyakovsky) inequality*

$$|\langle u, v \rangle| \leq \|u\| \cdot \|v\|,$$

perhaps one of the most important inequalities in mathematics, encountered everywhere from geometry to probability.

2. Cosine of the angle between two vectors $\cos \angle u, v = \frac{\langle u, v \rangle}{\|u\| \cdot \|v\|}$. The case of $\langle u, v \rangle = 0$ corresponds to a right angle between u, v , which are said to be *orthogonal* (sometimes denoted by $u \perp v$).
3. Not every norm defines an inner product! A norm satisfying the *parallelogram law*

$$2\|u\|^2 + 2\|v\|^2 = \|u + v\|^2 + \|u - v\|^2$$

defines an inner product through the *polarization identity*

$$\langle u, v \rangle = \frac{1}{4} (\|u + v\|^2 - \|u - v\|^2)$$

Theorem 1 (Pythagoras). Let $v_1, v_2, \dots, v_n \in V$ be pairwise orthogonal ($v_i \perp v_j$ for all $i \neq j$), i.e., $\langle v_i, v_j \rangle = \delta_{ij}$. Then,

$$\left\| \sum_{i=1}^n v_i \right\|^2 = \sum_{i=1}^n \|v_i\|^2.$$

For non-orthogonal vectors, $\left\| \sum_{i=1}^n v_i \right\|^2 \leq \sum_{i=1}^n \|v_i\|^2$ by virtue of triangle inequality.

4 Limits, convergence, Banach and Hilbert spaces

A sequence of vectors $v_1, v_2, \dots \in V$ is a *Cauchy sequence* if for any $\epsilon > 0 \exists N$ s.t.

$$\forall m, n > N \quad \|v_m - v_n\| < \epsilon,$$

i.e., $(v_n)_{n \geq 1}$ *converges in norm*: $\lim_{n \rightarrow \infty} \|v_m - v_n\| = 0$ (see below). This does not necessarily imply there exists a *limit* $\lim_{n \rightarrow \infty} v_n = v$ in V ; a classical example are sequences of rational numbers that can have an irrational limit. A *complete space* is a space in which every Cauchy sequence has a limit.

Convergence can refer to several different things:

1. *Convergence in norm* or *strong convergence* (sometimes denoted $v_n \xrightarrow{\|\cdot\|_V} v$ or simply $v_n \rightarrow v$) implies $\|v_n - v\| \xrightarrow{n \rightarrow \infty} 0$.
2. *Weak convergence* (sometimes denoted as $v_n \rightharpoonup v$) implies $\langle v_n, u \rangle \xrightarrow{n \rightarrow \infty} \langle v, u \rangle \quad \forall u \in V$.

Let $X \subseteq V$. The *closure* of X is defined by adding all the limit points of sequences in X , i.e.,

$$\overline{X} = X \cup \left\{ \lim_{n \rightarrow \infty} x_n, (x_n)_{n \geq 1} \subseteq X \right\}.$$

X is said to be *dense* in V if $\overline{X} = V$.

V is *separable* if it contains a countable dense subset.

Banach space is a complete normed vector space.

Hilbert space is a complete inner product space.

Remark. Sometimes Hilbert spaces are tacitly assumed separable, yielding the property of isometry to ℓ^2 .

5 Orthogonal bases

Let V be a Hilbert space and let $S \subseteq V$. The *span* of S , denoted $\text{span}(S) = \{\sum_{i=1}^n \alpha_i v_i : n \in \mathbb{N}, v_i \in S, \alpha_i \in \mathbb{C}\}$ is the set of all finite linear combinations from S .

S is *linearly independent* if its non-trivial finite linear combinations are always non-zero

$$\sum_{i=1}^n \alpha_i v_i \neq 0 \quad \forall n \in \mathbb{N}, v_i \in S, \alpha_i \in \mathbb{C} \text{ s.t. } \exists j : \alpha_j \neq 0.$$

S is *orthogonal* if $\langle u, v \rangle = 0 \quad \forall u, v \in S \text{ s.t. } u \neq v$.

S is *orthonormal* if it is orthogonal and in addition all vectors have unit length, i.e. $\|u\| = 1 \quad \forall u \in S$.

An *orthonormal basis* of V is a maximal orthonormal set.

Theorem 2. Every Hilbert space has an orthonormal basis¹ (not necessarily countable!)

¹This Theorem is a consequence of Zorn's Lemma, stating that if in a partially-ordered non-empty set X every totally-ordered subset $S \subset X$ has an upper bound in X , then $\exists \max(X)$.

Theorem 3. An orthonormal basis of Hilbert space V is countable iff V is separable.

Theorem 4. Let $S = \{v_1, v_2, \dots\} \subseteq V$ be an orthonormal subset of a separable Hilbert space V . Then, the following are equivalent:

1. S is maximal (i.e., a basis)
2. The set of finite linear combinations from S is dense in V .
3. $\sum_{i=1}^n \langle v, v_i \rangle v_i \xrightarrow{n \rightarrow \infty} v, \quad \forall v \in V$ (Fourier series converges in V)
4. $\sum_{i \geq 1} |\langle v, v_i \rangle|^2 = \|v\|^2, \quad \forall v \in V$ (Parseval identity)

We will discuss (3)-(4) in depth in the following.

6 Operators

An *operator* in this context is a map $A : U \rightarrow V$ between two spaces U and V (Banach or Hilbert), usually preserving some structure.

6.1 Operators between Banach spaces

Let $(U, \|\cdot\|_U)$ and $(V, \|\cdot\|_V)$ be Banach spaces and consider an operator $A : U \rightarrow V$.

A is *continuous* if it preserves convergence, i.e., $u_n \xrightarrow{\|\cdot\|_U} u \Rightarrow Au_n \xrightarrow{\|\cdot\|_V} Au$.

A is *bounded* if $\exists c > 0$ s.t. $\|Au\|_V \leq c\|u\|_U \quad \forall u \in U$.

A is *linear* if $A(\alpha u + \beta w) = \alpha Au + \beta Aw \quad \forall u, w \in U$ and $\alpha, \beta \in \mathbb{C}$.

A is an *isometry* if it is length-preserving, i.e. $\|Au\|_V = \|u\|_U$.

Theorem 5. For linear operators, boundedness and continuity are equivalent.

Remark. Pay attention that not every linear operator is bounded (and thus, not every linear operator is continuous). Consider for example $V = L^2([-\pi, +\pi])$ with the L_1 -norm and the derivative operator $A = \frac{d}{dx}$. Let $(f_n = e^{inx})_{n \geq 1}$. We have $\|f_n\| = 1$, however, $\|Af_n\| = \|inf_n\| = n \rightarrow \infty$, i.e., A is unbounded. Now construct $g_n = \frac{f_n}{\|Af_n\|} = \frac{1}{n}e^{inx}$. We have $g_n \rightarrow 0$ and at the same time $\|Ag_n\| = 1 \neq \|A0\| = 0$, i.e., A is not continuous (at zero vector) as it does not preserve limits.

6.2 Operators on Hilbert spaces

Let $(V, \langle \cdot, \cdot \rangle)$ be a Hilbert space and consider an operator $A : V \rightarrow V$.

A^* is *adjoint* to A if $\langle Au, v \rangle = \langle u, A^*v \rangle \quad \forall u, v \in V$.

A is *self-adjoint* if $A^* = A$, i.e. $\langle Au, v \rangle = \langle u, Av \rangle \quad \forall u, v \in V$.

A is *compact* if it maps weak limits to strong limits, i.e. $v_n \rightharpoonup v \Rightarrow Av_n \rightarrow Av$.

Remark. In the space of finite-dimensional real vectors, operators can be expressed as matrices: $\langle Au, v \rangle = (Au)^\top v = u^\top (A^\top v) = \langle u, A^\top v \rangle$. Adjoint for real matrices is the transpose, and a self-adjoint matrix is a symmetric matrix.

6.3 Functionals

6.3 Functionals

A *functional* is a map of the form $\phi : V \rightarrow \mathbb{C}$ on a Hilbert space V .

Dual (or *conjugate*) *space* to V is the space of linear continuous functionals on V , denoted

$$V^* = \{\phi : V \rightarrow \mathbb{C} \text{ linear+continuous}\}$$

Elements of V^* are called *dual vectors*.

Theorem 6 (Riesz-Fréchet). *Let $\phi \in V^*$. Then $\exists! u_\phi \in V$ s.t. $\phi(v) = \langle v, u_\phi \rangle \quad \forall v \in V$ (in other words: any dual vector can be represented as an inner product in V with a unique vector)*

Theorem 7. V^* is also a Hilbert space.

Proof outline. Riesz-Fréchet Theorem allows to define an operator $\Phi : V \rightarrow V^*$ associating a vector $u \in V$ with a linear functional $\Phi(u)$ acting on $v \in V$ as $\Phi(u)v = \langle v, u \rangle$. This functional is bijective (due to uniqueness of representation) and conjugate linear,

$$\Phi(\alpha v + \beta w) = \overline{\alpha}\Phi(v) + \overline{\beta}\Phi(w).$$

The inner product on V^* is pulled back from V by Φ^{-1} (which exists due to bijectivity),

$$\langle \phi, \psi \rangle_{V^*} = \langle \Phi^{-1}\phi, \Phi^{-1}\psi \rangle_V, \quad \forall \phi, \psi \in V^*.$$

It also appears that Φ is an isometry between V and V^* . □

Remark. *The familiar concept of vector coordinates is in fact correctly expressed in the form of dual vectors, i.e., functionals acting on vectors and producing scalars.*

Remark. *The delta function is defined as the linear functional $\delta f = f(0)$, and can be represented as inner product with an impulse*

$$\delta f = \langle \delta, f \rangle, \quad \delta(x) = \begin{cases} \infty & x = 0 \\ 0 & x \neq 0 \end{cases}$$

7 Eigenfunctions and eigenvalues

Let $A : V \rightarrow V$ an operator on Hilbert space V . A vector $v \neq 0$ satisfying for some λ

$$Av = \lambda v$$

is called an *eigenvector* (or *eigenfunction* in case V is a space of functions) of A , and λ is the corresponding *eigenvalue*.

Remark. *Note that eigenvectors are defined up to scale: if v is an eigenvector of A , so is αv for any $\alpha \neq 0$, since we can multiply both sides of the equation by $A\alpha v = \lambda\alpha v$ by α . It is common to assume eigenvectors of unit length, i.e. $\|v\| = 1$.*

Theorem 8. Self-adjoint operators have real eigenvalues.

Proof. Let $Av = \lambda v$, with $v \neq 0$. Then, since $A = A^*$,

$$\begin{aligned}\langle Av, v \rangle &= \langle v, Av \rangle \\ \langle \lambda v, v \rangle &= \langle v, \lambda v \rangle \\ \lambda \langle v, v \rangle &= \bar{\lambda} \langle v, v \rangle\end{aligned}$$

Since $v \neq 0$, we can divide both sides by $\|v\|^2 = \langle v, v \rangle > 0$ to get $\lambda = \bar{\lambda}$, from which it follows that $\lambda \in \mathbb{R}$. □

Theorem 9. *Eigenfunctions of self-adjoint operators corresponding to different eigenvalues are orthogonal.*

Proof. Let $Av = \lambda v$ and $Aw = \mu w$ with $\lambda \neq \mu$ and $v, w \neq 0$. Then, since $A = A^*$,

$$\begin{aligned}\langle Av, w \rangle &= \langle v, Aw \rangle \\ \langle \lambda v, w \rangle &= \langle v, \mu w \rangle\end{aligned}$$

Since, by Theorem 8, the eigenvalues λ, μ are real, we can take the scalars outside the inner product without conjugation,

$$\begin{aligned}\lambda \langle v, w \rangle &= \mu \langle v, w \rangle \\ (\lambda - \mu) \langle v, w \rangle &= 0\end{aligned}$$

Since $\lambda \neq \mu$, it must be that $\langle v, w \rangle = 0$, i.e., $v \perp w$. □

Example Let $L^2([-\pi, +\pi]) = \left\{ f : [-\pi, +\pi] \rightarrow \mathbb{C} \text{ s.t. } \int_{-\pi}^{+\pi} |f(x)|^2 dx < \infty \right\}$ be the space of square-integrable periodic functions with standard inner product $\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{+\pi} f(x) \overline{g(x)} dx$. Consider the Laplacian operator (second-order derivative) $\Delta = -\frac{d^2}{dx^2}$. Note that we define the Laplacian with negative sign here, for reasons that will become apparent in the following.

First, verify that Δ is self-adjoint. From product differentiation rule (or integration by parts) and periodic boundary conditions $f(-\pi) = f(+\pi)$, we have

$$\begin{aligned}\int_{-\pi}^{+\pi} (fg)'(x) dx &= \int_{-\pi}^{+\pi} f'(x)g(x) dx + \int_{-\pi}^{+\pi} f(x)g'(x) dx = f(x)g(x) \Big|_{-\pi}^{+\pi} \stackrel{\text{period.}}{=} 0 \\ \int_{-\pi}^{+\pi} f(x)g'(x) dx &= - \int_{-\pi}^{+\pi} f'(x)g(x) dx.\end{aligned}$$

(for simplicity, we ignore complex conjugates). Applying this result to $f'g'$ we have

$$- \int_{-\pi}^{+\pi} f(x)g''(x) dx = \int_{-\pi}^{+\pi} f'(x)g'(x) dx = - \int_{-\pi}^{+\pi} f''(x)g(x) dx$$

from which self-adjointness $\langle \Delta f, g \rangle = \langle f, \Delta g \rangle$ follows.

Let us now look at the eigenfunctions. It is easy to verify that $\Delta e^{inx} = n^2 e^{inx}$ (here $i = \sqrt{-1}$), i.e., the eigenfunctions have the form e^{inx} with corresponding real eigenvalues n^2 , for $n = 0, 1, \dots$. To verify

orthogonality, write

$$\begin{aligned}\langle e^{inx}, e^{imx} \rangle &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} e^{inx} e^{-imx} dx \\ &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} e^{i(n-m)x} dx = \frac{1}{2\pi} 2\pi \delta_{mn} = \delta_{mn}.\end{aligned}$$

□

An operator $A : V \rightarrow V$ is compact iff it can be written in the form

$$Aw = \sum_{n \geq 1} \sigma_n \langle v_n, w \rangle u_n, \quad \forall w \in V$$

(this is an alternative definition of compactness equivalent to the one based on convergence we have already seen; we leave this equivalence without proof). $\{\sigma_n\}_{n \geq 1}$ are the *singular values* and $\{v_n\}_{n \geq 1}, \{u_n\}_{n \geq 1}$ are the corresponding (left- and right-) *singular vectors* of A .

Remark. Singular values $\{\sigma_n\}_{n \geq 1}$ can accumulate only at zero, i.e., starting from some N , $\sigma_{n \geq N} = 0$ (in this case $\text{rank}(A) = N$).

Remark. An $m \times n$ matrix A can be written in the form

$$A = U \Sigma V^* = \begin{pmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{pmatrix} \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{pmatrix} \begin{pmatrix} - & \bar{v}_1 & - \\ & \vdots & \\ - & \bar{v}_n & - \end{pmatrix},$$

where U, Σ , and V are matrices of sizes $m \times n$, $n \times n$, and $n \times n$.

Remark. If A is in addition self-adjoint, it can be written as

$$Aw = \sum_{n \geq 1} \lambda_n \langle v_n, w \rangle v_n, \quad \forall w \in V.$$

For matrices, this corresponds to the well-known fact that symmetric matrices have orthogonal eigendecomposition.

Theorem 10 (Spectral Theorem). A compact self-adjoint operator $A : V \rightarrow V$ has eigenvectors $\{v_\lambda\}$ with corresponding eigenvalues $\{\lambda\}$ satisfying $Av_\lambda = \lambda v_\lambda$, which form an orthonormal basis of V . Furthermore, the basis is discrete.

8 Fourier analysis

Let $\{v_\alpha\}$ be an orthonormal basis in V . Then, $u \in V$ can be expressed as a *Fourier series*

$$u = \sum_{\alpha} \langle u, v_\alpha \rangle v_\alpha$$

The coefficients $\langle u, v_\alpha \rangle = \hat{u}_\alpha$ in the above series are called *Fourier coefficients* (or *transform*) of u .

8.1 Fourier transform

Remark

1. Usually, the term ‘Fourier series’ refers to trigonometric basis (sine, cosine, or complex exponential, as we will see in the following). However, the concept is general and applies to any orthonormal basis.
2. For vectors, the Fourier decomposition can be expressed as

$$u = \begin{pmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{pmatrix} \begin{pmatrix} - & \bar{v}_1 & - \\ & \vdots & \\ - & \bar{v}_n & - \end{pmatrix} u$$

from which it is evident that it is a unitary operation (see below).

Theorem 11 (Parseval’s identity). *Let $u = \sum_{\alpha} \hat{u}_{\alpha} v_{\alpha}$ and $w = \sum_{\alpha} \hat{w}_{\alpha} v_{\alpha}$ be Fourier series of $u, w \in V$ w.r.t. the orthonormal basis $\{v_{\alpha}\}$. Then $\langle u, w \rangle = \sum_{\alpha} \hat{u}_{\alpha} \bar{\hat{w}}_{\alpha}$.*

In other words, we can define a map $V \ni u \mapsto \hat{u} = \{\langle u, v_{\alpha} \rangle\} \in \ell^2$ from vectors to (square summable) sequences. This map is an isometry:

$$\|u\|_V^2 = \sum_{\alpha} |\langle u, v_{\alpha} \rangle|^2 = \sum_{\alpha} |\hat{u}_{\alpha}|^2 = \|\hat{u}\|_{\ell^2}^2.$$

This, in turn, is nothing else but the application of the Pythagorean theorem (possibly in infinite dimension),

$$\|u\|^2 = \left\| \sum_{\alpha} \langle u, v_{\alpha} \rangle v_{\alpha} \right\|^2 = \sum_{\alpha} \|\langle u, v_{\alpha} \rangle v_{\alpha}\|^2 = \sum_{\alpha} |\langle u, v_{\alpha} \rangle|^2,$$

where we used the orthonormality of the basis $\{v_{\alpha}\}$.

Example Let $f \in L^2([-\pi, +\pi])$ be a square-integrable function. We assume the standard inner product $\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{+\pi} f(x) \overline{g(x)} dx$ and the basis $\{e^{inx}\}_{n \geq 1}$. Then,

$$f(x) = \sum_{n \geq 1} \frac{1}{2\pi} \int_{-\pi}^{+\pi} f(y) e^{-iny} dy e^{inx},$$

which is the classical statement of the celebrated Fourier’s result that any function can be expressed as a linear combination of sinusoids.

8.1 Fourier transform

Let $f \in L^2(\mathbb{R})$ be square-integrable function. Its *Fourier transform* is given by

$$\hat{f}(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-ix\omega} dx,$$

and the *inverse Fourier transform*

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

Remark. *There exist several definitions of the Fourier transform, typically differing up to normalization and scale of the frequency. We chose our definition to make the Fourier transform an isometry (see below).*

Considered as an operator on the space of square integrable functions $\mathcal{F} : L^2(\mathbb{R}) \rightarrow L^2(\mathbb{R})$, the Fourier transform is an isometry, a fact captured by the Parseval identity, which assumes the following form

$$\int_{-\infty}^{\infty} |f(x)|^2 dx = \int_{-\infty}^{\infty} |\hat{f}(\omega)|^2 d\omega$$

in the continuous case.

Table 1 summarizes some important properties of the Fourier transform.

	Spatial domain	Frequency domain
	$f(x)$	$\hat{f}(\omega)$
Linearity	$\alpha f(x) + \beta g(x)$	$\alpha \hat{f}(\omega) + \beta \hat{g}(\omega)$
Shift	$f(x \pm x_0)$	$e^{\pm i x_0 \omega} \hat{f}(\omega)$
Modulation	$e^{\pm i \omega_0 x} f(x)$	$\hat{f}(\omega \mp \omega_0)$
Delta function	$\delta(x)$	1
Convolution	$(f \star g)(x)$	$\hat{f}(\omega) \cdot \hat{g}(\omega)$
Derivative	$\frac{d}{dx} f(x)$	$i\omega \hat{f}(\omega)$
Moment	$x f(x)$	$i \frac{d}{d\omega} \hat{f}(\omega)$

Table 1: Some properties of the Fourier transform.

9 Heat equation

Consider the following partial differential equation (Dirichlet problem) called the *heat equation*

$$\begin{cases} \partial_t f(x, t) = -c \Delta f(x, t) \\ f(x, 0) = g(x) \end{cases} \quad (\text{initial conditions})$$

on a circle, where $f : S^1 \times [0, \infty) \rightarrow \mathbb{R}$ (periodic in the first coordinate) represents the temperature at point x in time t , $\Delta = -\frac{\partial^2}{\partial x^2}$ is the (negative) one-dimensional Laplacian operator, and $g(x)$ is the initial temperature distribution at time $t = 0$ (there is no boundary on S^1 , hence no boundary conditions). In the following analysis, we assume for simplicity $c = 1$.

Remark. The heat equation is a differentiable form of Newton's law of cooling, stating that the rate of change of the temperature of an object is proportional to the difference between its own temperature and the temperature of the surrounding. In the formula $f_t = -c \Delta f$, the rate of change is captured by the temporal derivative f_t , while the Laplacian $-\Delta f$ is the difference between the temperature of a point and its infinitesimal surrounding. The proportion coefficient c is called the thermal diffusivity constant. As a sanity check, we can verify the physical units in this equation: $f_t [1/\text{sec}] = -c [m^2/\text{sec}] \Delta f [1/m^2]$.

Fourier analysis has been originally developed for the solution of this kind of PDEs, and we will exemplify how it is used. First, assume the solution has a separable form $f(x, t) = X(x)T(t)$ into a spatial part X and

temporal part T . Then (assuming X, T never vanish), since

$$\partial_t f + \Delta f = X(\partial_t T) + (\Delta X)T = 0,$$

holds for any x, t , it follows that

$$\frac{\Delta X}{X} = -\frac{\partial_t T}{T} = \lambda \text{ (some constant).}$$

In other words, the spatial and temporal part of the solution are eigenfunctions of the Laplacian and first-order derivative operators, respectively:

$$\Delta X = \lambda X \quad \partial_t T = -\lambda T,$$

which we can express as

$$\Delta e^{inx} = n^2 e^{inx} \quad \partial_t e^{-n^2 t} = -n^2 e^{-n^2 t},$$

and $\lambda = n^2$, $n = 0, 1, \dots$ is the corresponding eigenvalue. Note that the eigenvalues of the Laplacian are nonnegative, i.e. Δ is a positive-semidefinite operator – this is the main reason for defining it with negative sign.

Hence, solutions of the equation are of the form $f_n(x, t) = e^{inx} e^{-n^2 t}$. However, due to linearity, $f_m + f_n$ is also a solution, leading to the general expression

$$f(x, t) = \sum_{n \geq 0} a_n e^{inx} e^{-n^2 t}.$$

In order to find a unique solution, we must use the initial condition. Note that since $\{e^{inx}\}_{n \geq 1}$ is an orthonormal basis on $L^2(S^1)$, we can express the initial condition as a Fourier series

$$g(x) = \sum_{n \geq 0} \langle g, e^{inx} \rangle e^{inx}.$$

Since $f(x, 0) = g(x)$ we can identify $a_n = \hat{g}_n = \langle g, e^{inx} \rangle$ in the above equations, yielding

$$\begin{aligned} f(x, t) &= \sum_{n \geq 0} \frac{1}{2\pi} \int_{-\pi}^{+\pi} g(y) e^{-iny} dy e^{inx} e^{-n^2 t} \\ &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} g(y) \underbrace{\sum_{n \geq 0} e^{-n^2 t} e^{-in(x-y)}}_{h_t(x-y)} dy = g \star h_t, \end{aligned}$$

where the exchange of summation and integration is done under the usual Fubini-Tonelli conditions.

Remark. h_t is called the *fundamental solution of the heat equation* or the *heat kernel*. In particular, for $g(x) = \delta(x)$ (impulse initial condition), from the definition of the delta function $\langle f, \delta \rangle = f(0)$, we get $\langle \delta, e^{inx} \rangle = e^{in0} = 1$, i.e. $a_{n \geq 0} = 1$, which implies

$$f(x, t) = \sum_{n \geq 0} e^{-n^2 t} e^{inx} = h_t(x).$$

In signal processing terms, h_t is thus the impulse response.

The operator $H^t : g(x) \mapsto f(x, t) = (g \star h_t)(x)$ is called the *heat operator*. It is easy to see that the

eigenfunctions of the heat operator are the same as those of the Laplacian,

$$\begin{aligned} H^t e^{imx} &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} e^{imy} \sum_{n \geq 0} e^{-n^2 t} e^{in(x-y)} dy \\ &= \frac{1}{2\pi} \sum_{n \geq 0} e^{-n^2 t} \underbrace{\int_{-\pi}^{+\pi} e^{i(m-n)y} dy}_{2\pi \delta_{mn}} e^{inx} = e^{-m^2 t} e^{imx}, \end{aligned}$$

while the eigenvalues are $e^{-m^2 t}$. Thus, the heat operator is given as the exponential of the Laplacian $H^t = e^{-t\Delta}$, and is compact and self-adjoint.

Remark. Applying functions to operators is understood as applying functions to their eigenvalues. Let $A : V \rightarrow V$ self-adjoint compact operator and $f : \mathbb{R} \rightarrow \mathbb{R}$ some function. Since we can express A in terms of its orthonormal eigenvectors $\{v_n\}_{n \geq 1}$, we have

$$f(A)u = \sum_{n \geq 1} f(\lambda_n) \langle v_n, u \rangle v_n.$$

The function $f(\lambda)$ can be thus given the interpretation of a spectral filter operating on the eigenvalues $\{\lambda_n\}_{n \geq 1}$ of A . We will exploit this intuition extensively in the following.

Remark. The heat operator $H^t = e^{-t\Delta}$ can be interpreted as a low-pass filter, where the attenuation becomes more significant as t grows, since Δ is positive-semidefinite. This corresponds to the intuitive understanding of diffusion as a ‘blurring’ process.

10 Wave equation

Let us now look at the wave equation,

$$f_{tt}(x, t) = -v^2 \Delta f(x, t)$$

where $f(x, t)$ denotes the wave function, and Δ is the (negative) second derivative as before.

Remark. The wave propagation is modeled as a displacement of masses connected between each other by a system of springs. The wave equation encodes a differential version of Newton’s second law: the right hand is the acceleration proportional to the net force, in this case elastic force of the spring given by Hooke’s law. The proportion coefficient v is interpreted as the wave speed and has the units of $[m/sec]$; in the following we will assume for simplicity $v = 1$.

Similarly to the heat equation, the wave equation is solved by separation of variables, assuming solution of the form $f(x, t) = X(x)T(t)$, which results in

$$\frac{\partial_{tt} T}{T} = -\frac{\Delta X}{X} = -\lambda.$$

The equation $\Delta X(x) = \lambda X(x)$ satisfied by the spatial part is called the *Helmholtz equation*; its solutions are the Laplacian eigenfunctions that can be interpreted as ‘vibration modes’ of the domain, and corresponding eigenvalues are the vibration frequencies.

11 Vector calculus

11.1 Scalar fields

We are interested in analyzing functions on various domains (Euclidean space, manifolds, and graphs), and can define various classes of functions satisfying some regularity properties. For the sake of simplicity, we avoid here a formal definition of distributions, weak derivatives, and Sobolev spaces, and assume our functions are ‘good enough’, in particular, smooth (infinitely continuously differentiable) and square-integrable. We denote by $\mathcal{C}^\infty(\mathbb{R}^n)$ the set of smooth real functions on \mathbb{R}^n , to which we also refer as *scalar fields*. We use the standard inner product

$$\langle f, g \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)} = \int_{\mathbb{R}^n} f(x)g(x) dx,$$

(here dx denotes the volume element on \mathbb{R}^n), and the respective norm $\|f\|_{\mathcal{C}^\infty(\mathbb{R}^n)} = \langle f, f \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)}^{1/2}$.

11.2 Vector fields

Let $x \in \mathbb{R}^n$. The set $T_x\mathbb{R}^n = \{(x, v) : v \in \mathbb{R}^n\}$ is called the *tangent space of \mathbb{R}^n at x* (the reason for this term will become later when we discuss manifolds). Whenever no confusion arises, we will use the short notation $v \in T_x\mathbb{R}^n$ referring to (x, v) .

$T_x\mathbb{R}^n$ can be easily made into a vector space by defining vector addition operation $(x, v) + (x, w) = (x, v+w)$ (note that $(x, v) + (y, w)$ makes no sense: only endpoints $v + w$ but not the origins $x + y$ can be added) and vector by scalar multiplication $\alpha \cdot (x, v) = (x, \alpha \cdot v)$. Elements of the tangent space are called *tangent vectors*. An inner product between tangent vectors is defined simply as $\langle v, w \rangle_{T_x\mathbb{R}^n} = \langle v, w \rangle_{\mathbb{R}^n}$.

A *vector field* is a function F assigning each point to a tangent vector, $\mathbb{R}^n \ni x \mapsto F(x) \in T_x\mathbb{R}^n$. Let $\{e_{1,x}, \dots, e_{n,x}\}$ be the standard basis on $T_x\mathbb{R}^n$ (while not the case for Euclidean space, we insist in this notation that the basis is dependent on x , which is the case for manifolds). We can then express F as

$$F(x) = F_1(x)e_{1,x} + \dots + F_n(x)e_{n,x}.$$

F is said to be continuous, differentiable, smooth, etc. if all of its components F_1, \dots, F_n are such; in the following we will tacitly assume that all vector fields are \mathcal{C}^∞ .

The disjoint union $T\mathbb{R}^n = \bigsqcup_{x \in \mathbb{R}^n} T_x\mathbb{R}^n$ is called the *tangent bundle* of \mathbb{R}^n . A vector field can thus be seen as a map $F : \mathbb{R}^n \rightarrow T\mathbb{R}^n$ between the space and its tangent bundle; we denote by $\mathcal{C}^\infty(T\mathbb{R}^n)$ the set of all smooth vector fields on \mathbb{R}^n , modulo constants. This is also a vector space, which can be equipped with the inner product

$$\langle F, G \rangle_{\mathcal{C}^\infty(T\mathbb{R}^n)} = \int_{\mathbb{R}^n} \langle F(x), G(x) \rangle_{T_x\mathbb{R}^n} dx,$$

and the respective norm $\|F\|_{\mathcal{C}^\infty(T\mathbb{R}^n)} = \langle F, F \rangle_{\mathcal{C}^\infty(T\mathbb{R}^n)}^{1/2}$.

11.3 Gradient

Let $f \in \mathcal{C}^\infty(\mathbb{R}^n)$ be a scalar field. In order to determine how f changes at a point, the usual notion of derivative suggests to look at the *differential* $df(x) = f(x + dx) - f(x)$, where dx is an infinitesimal step. We can redefine this notion in terms of tangent vectors (it would not be possible otherwise on manifolds), as follows: the

11.4 Divergence

differential $df(x) : T_x\mathbb{R}^n \rightarrow \mathbb{R}$ of f at x is a linear functional acting on vector fields as

$$df(x)v = \langle \nabla f(x), v \rangle_{T_x\mathbb{R}^n}$$

for any $v \in T_x\mathbb{R}^n$. Thinking of this functional as a dual vector (element of the *cotangent space* $(T_x\mathbb{R}^n)^* = T_x^*\mathbb{R}^n$), its representation in the Riezs-Fréchet sense is the tangent vector $\nabla f(x)$, called the *gradient* of f at x . The differential can be thought of as the map $df : \mathcal{C}^\infty(\mathcal{X}) \rightarrow T^*\mathcal{X}$ to the cotangent bundle, sometimes referred to as the *differential map*.

Remark. We stress that vectors should be correctly treated as abstract objects rather than their coordinates in some basis, and this will be especially important for manifolds. Given $\{e_{1,x}, \dots, e_{n,x}\}$ the standard basis on $T_x\mathbb{R}^n$, one can express the gradient by applying $\langle \nabla f(x), e_{i,x} \rangle_{T_x\mathbb{R}^n} = \frac{\partial}{\partial x_i} f(x)$. This leads to the usual way of thinking of the gradient as a vector of partial derivatives,

$$\nabla f(x) = \left(\frac{\partial}{\partial x_1} f(x), \dots, \frac{\partial}{\partial x_n} f(x) \right).$$

Remark. Given a scalar field $f \in \mathcal{C}^\infty(\mathbb{R}^n)$ and a vector field $F \in \mathcal{C}^\infty(T\mathbb{R}^n)$, we can define a map $\nabla_F f$, associating with $f(x)$ its derivative $df(x)F(x)$ in the direction of the tangent vector $F(x)$,

$$(\nabla_F f)(x) = \langle \nabla f(x), F(x) \rangle_{T_x\mathbb{R}^n}, \quad \forall x \in \mathbb{R}^n.$$

$\nabla_F f$ is called the covariant derivative of f along F and is a generalization of the traditional notion of directional derivative.

The *gradient operator* is a map $\nabla : \mathcal{C}^\infty(\mathbb{R}^n) \rightarrow \mathcal{C}^\infty(T\mathbb{R}^n)$ associating $f(x)$ with a tangent vector $\nabla f(x)$ that indicates the direction of the steepest change of f at x .

11.4 Divergence

It is convenient to think of a vector field $F \in \mathcal{C}^\infty(T\mathbb{R}^n)$ as of a *flow* of particles in space, where $F(x)$ indicates in which direction and at what speed a particle at point x moves. The net flow of a field F at a point x is called the *divergence* of F and denoted $\operatorname{div} F(x)$. The sign of the divergence allows to distinguish between field ‘sources’ and ‘sinks’, i.e., if a field is created or disappears.

Remark. In the standard basis, the divergence is expressed as

$$\operatorname{div} F(x) = \frac{\partial}{\partial x_1} F_1(x) + \dots + \frac{\partial}{\partial x_n} F_n(x),$$

which is often symbolically represented as $\langle F(x), \nabla \rangle_{T_x\mathbb{R}^n}$ or $\nabla \cdot F$.

The (negative) divergence operator $-\operatorname{div} : \mathcal{C}^\infty(T\mathbb{R}^n) \rightarrow \mathcal{C}^\infty(\mathbb{R}^n)$ is the adjoint of the gradient operator, in the following sense

$$\langle F, \nabla f \rangle_{\mathcal{C}^\infty(T\mathbb{R}^n)} = \langle -\operatorname{div} F, f \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)}$$

(note that the inner products in the left- and right-hand sides of the equations are taken on vector and scalar fields, respectively).

11.5 Laplacian

Theorem 12 (Gauss-(Ostrogradsky-Stokes) or simply Divergence theorem). *Let $\Omega \subseteq \mathbb{R}^n$ be a region in space with boundary $\partial\Omega$. Then,*

$$\int_{\Omega} \operatorname{div} F(x) dx = \int_{\partial\Omega} \langle F(x), \hat{n}(x) \rangle dx,$$

where $\hat{n}(x)$ denotes the unit normal vector to the boundary surface $\partial\Omega$ at point x on thereon, and with some of abuse of notation dx in the left- and right-hand sides of the equation denotes n - and $(n-1)$ -dimensional volume elements, respectively.

The divergence theorem is a mathematical statement of the physical *conservation law* that, in the absence of the creation or destruction of matter, the density within a region of space can change only by having it flow into or away from the region through its boundary.

11.5 Laplacian

We define the *Laplacian* of f as

$$\Delta f = -\operatorname{div}(\nabla f)$$

The operator $\Delta : \mathcal{C}^\infty(\mathbb{R}^n) \rightarrow \mathcal{C}^\infty(\mathbb{R}^n)$ is self-adjoint (we have already seen this in 1D case), which is easy to verify from the above definition,

$$\langle \Delta f, g \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)} = \langle -\operatorname{div}(\nabla f), g \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)} = \langle \nabla f, \nabla g \rangle_{\mathcal{C}^\infty(T\mathbb{R}^n)} = \langle f, \Delta g \rangle_{\mathcal{C}^\infty(\mathbb{R}^n)}.$$

W.r.t. the standard basis, the Laplacian can be familiarly written as the (negative) sum of second-order derivatives,

$$\nabla f(x) = - \left(\frac{\partial^2}{\partial x_1^2} f(x) + \dots + \frac{\partial^2}{\partial x_n^2} f(x) \right).$$

The Laplacian arises ubiquitously in mathematical physics equations and is prominently featured in this course; it therefore deserves devoting some attention trying to better understand the intuition behind it. Geometrically, the $\Delta f(x)$ can be interpreted as the difference between $f(x)$ and the local average of f in an infinitesimal neighborhood (sphere of radius $r \rightarrow 0$) of x . We can formalize it as the following

Theorem 13. *Let $S_r^{n-1}(x)$ be an $(n-1)$ -dimensional sphere of radius r around x . Then*

$$\Delta f(x) = -\frac{2n}{r^2} \frac{1}{\operatorname{vol}(S_r^{n-1})} \int_{S_r^{n-1}(x)} (f(y) - f(x)) dy + \mathcal{O}(r^{n+2}).$$

Proof. For simplicity and w.l.o.g., we will prove the case for $x = 0$. From second-order Taylor expansion we have explicitly in coordinates

$$f(x) = f(0) + \underbrace{\sum_{i=1}^n g_i x_i}_{x^\top \nabla f(0)} + \frac{1}{2} \underbrace{\sum_{i,j=1}^n h_{ij} x_i x_j}_{x^\top \nabla^2 f(0) x} + \mathcal{O}(\|x\|^3),$$

where $g = \nabla f(x)$ and $H = \nabla^2 f(0)$ are the gradient and Hessian of f at 0, respectively.

Let us now integrate $f(x) - f(0)$ on $S_r^{n-1} = S_r^{n-1}(0)$. It is apparent immediately, due to symmetry consid-

11.5 Laplacian

erations, that odd terms of the form $\int_{S_r^{n-1}} g_i x_i dx$ and $\int_{S_r^{n-1}} h_{ij} x_i x_j dx$ for $i \neq j$ vanish, while

$$\int_{S_r^{n-1}} x_i^2 dx = \frac{1}{n} \int_{S_r^{n-1}} \sum_{i=1}^n x_i^2 dx = \frac{1}{n} \int_{S_r^{n-1}} \|x\|^2 dx = \frac{r^2}{n} \text{vol}(S_r^{n-1}).$$

This leaves us with

$$\begin{aligned} \int_{S_r^{n-1}} (f(x) - f(0)) dx &= \frac{1}{2} \sum_{i=1}^n h_{ii} \frac{r^2}{n} \text{vol}(S_r^{n-1}) + \mathcal{O}(r^{(n-1)+3}) \\ &= -\frac{r^2}{2n} \text{vol}(S_r^{n-1}) \Delta f(0) + \mathcal{O}(r^{n+2}), \end{aligned}$$

from which the desired result follows. \square

From this geometric interpretation of the Laplacian as the difference between the value of a function at a point and the average value in an infinitesimal neighborhood, the following well-known fact immediately follows:

Theorem 14. *The (Euclidean) Laplacian is rotation-invariant.*

Proof. We need to prove invariance under a change of coordinates Ax , where $AA^\top = A^\top A = I$ (orthogonal matrix describing rotation). To this end, we write the Laplacian as the (negative) trace of the Hessian, $\Delta f(x) = -\text{trace}(\nabla^2 f(x))$ (when representing the Hessian as a matrix w.r.t. the standard basis, its diagonal contains the second order derivatives $\frac{\partial^2}{\partial x_i^2} f(x)$). Then, applying the chain rule, we have

$$\begin{aligned} \nabla_x f(Ax) &= A^\top \nabla_{Ax} f(Ax) \\ \nabla_x^2 f(Ax) &= A^\top \nabla_{Ax}^2 f(Ax) A. \end{aligned}$$

Using matrix commutativity under trace we get

$$\begin{aligned} \Delta_x f(Ax) &= -\text{trace}(A^\top \nabla_{Ax}^2 f(Ax) A) \\ &= -\text{trace}(\nabla_{Ax}^2 f(Ax) A A^\top) \\ &= -\text{trace}(\nabla_{Ax}^2 f(Ax)) = \Delta_{Ax} f(Ax). \end{aligned}$$

\square

Laplacian eigenfunctions. The quadratic functional

$$\mathcal{E}(f) = \|\nabla f\|_{C^\infty(T\mathbb{R}^n)}^2 = \int_{\mathbb{R}^n} \|\nabla f(x)\|_{T_x \mathbb{R}^n}^2 dx \geq 0$$

is called the *Dirichlet energy* in physics and quantifies how much f varies ($f = \text{const}$ has $\mathcal{E} = 0$). The eigenfunctions of the Laplacian can be obtained as the minimizers of the Dirichlet energy with orthogonality constraint,

$$\phi_k = \arg \min_{\phi} \|\nabla \phi\|_{C^\infty(T\mathbb{R}^n)}^2 \quad \text{s.t.} \quad \phi_k \perp \phi_{k-1}, \dots, \phi_1, \quad \text{and} \quad \|\phi_k\|_{C^\infty(\mathbb{R}^n)} = 1,$$

which allows to interpret it as the smoothest orthonormal basis. We will see in what follows how such functions look like exactly.

12 Convolution

Many textbooks define convolutions and Fourier transforms axiomatically, without providing a motivation where these constructions come from. Here, we would like to take a different avenue, defining shift-equivariance from which the convolution operation and the Fourier basis emerge by themselves. For simplicity, we analyze the discrete case; a continuous analogy can be derived in a similar manner, though somewhat more elaborate.

12.1 Circulant matrices

Let $x = (x_0, \dots, x_{n-1})^\top$, $n \in \mathbb{Z}_n$ be a vector with indices defined mod n , modeling discrete signals on the circle. We define a *circulant matrix*

$$C(x) = \begin{pmatrix} x_0 & x_{n-1} & \dots & x_1 \\ x_1 & x_0 & & x_2 \\ \vdots & & \ddots & \vdots \\ x_{n-1} & x_{n-2} & \dots & x_0 \end{pmatrix}$$

comprising circularly shifted versions of vector x as its columns. Elements of $C(x)$ have the form $c_{ij}(x) = x_{i-j \bmod n}$. Circulant matrices describe *circular convolution* operations $x \circledast y = C(x)y$, given by

$$(x \circledast y)_i = \sum_{k=0}^{n-1} c_{ik}(x)y_k = \sum_{k=0}^{n-1} x_{i-k \bmod n} y_k$$

It is easy to see that circular convolution is

1. *Linear*: $(\alpha x + \beta y) \circledast z = \alpha x \circledast z + \beta y \circledast z$
2. *Commutative*: $x \circledast y = y \circledast x$
3. *Associative*: $(x \circledast y) \circledast z = x \circledast (y \circledast z)$

Theorem 15. *Product of circulant matrices is also a circulant matrix.*

Proof. follows immediately from the above properties:

$$C(x \circledast y)z = (x \circledast y) \circledast z \stackrel{\text{assoc.}}{=} x \circledast (y \circledast z) = C(x)C(y)z$$

□

12.2 Shift operators

A particular kind of circulant matrix are *shift operators*. We define the *right shift* S and *left shift* S^\top as

$$S = \begin{pmatrix} 0 & \dots & 0 & 1 \\ 1 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 1 & 0 \end{pmatrix} \quad S^\top = \begin{pmatrix} 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ 1 & 0 & \dots & 0 \end{pmatrix}$$

12.2 Shift operators

Elements of S have the form $s_{ij} = \delta_{i-j-1 \bmod n}$, and its application on x has the form $(Sx)_i = x_{i-1 \bmod n}$. Similarly, $(S^\top x)_i = x_{i+1 \bmod n}$.

Shift operator is orthogonal $S^\top S = SS^\top = I$, representing the intuition that if we first shift left and shift right the results remains the same. Consequently, we expect S to have orthonormal eigenvectors, as we will see in the following.

Theorem 16 (Circulant matrices commute with shift). *A is circulant iff it commutes with S , i.e.*

$$AS = SA.$$

Proof. \Rightarrow Assume $A = C(x)$ is circulant, for some x . Then

$$\begin{aligned} (SC(x))_{ij} &= \sum_k s_{ik} c_{kj}(x) = \sum_k \delta_{i-k-1 \bmod n} x_{k-j \bmod n} = x_{i-j-1 \bmod n} \\ (C(x)S)_{ij} &= \sum_k c_{ik}(x) s_{kj} = \sum_k x_{i-k \bmod n} \delta_{k-j-1 \bmod n} = x_{i-j-1 \bmod n}, \end{aligned}$$

from where the desired commutativity $SC(x) = C(x)S$ follows.

\Leftarrow Assume $AS = SA$ for some A . Then

$$\begin{aligned} (SA)_{ij} &= \sum_k s_{ik} a_{kj} = \sum_k \delta_{i-k-1 \bmod n} a_{kj} = a_{i-1 \bmod n, j} \\ (AS)_{ij} &= \sum_k a_{ik} s_{kj} = \sum_k a_{ik} \delta_{k-j-1 \bmod n} = a_{i, j+1 \bmod n}. \end{aligned}$$

The equality $a_{i-1 \bmod n, j} = a_{i, j+1 \bmod n}$ implies A is circulant. \square

Theorem 17. *Convolution is a linear shift equivariant operation.*

Proof. Obtained immediately using associativity and commutativity properties:

$$(Sx) \circledast y = y \circledast (Sx) = C(y)Sx = SC(y)x = S(x \circledast y),$$

i.e. convolving a shifted version of x with y is equivalent to shifting the convolution of x and y . \square

Remark. *The aforementioned properties is often incorrectly referred to as shift-invariance, especially in the signal processing literature. We emphasize the difference between the two terms. Shift invariance implies $CSx = Cx$, i.e. the result of the operator C is unaffected by the shift of the input x (hence the name, ‘invariant’ means ‘unchanged’). Shift equivariance implies that the output of the operator C changes the same way as the input, $CSx = SCx$ (hence the name, ‘equivariant’ means ‘changes in the same way’).*

Note that the two above theorems give us a way to define the convolution as a linear shift-equivariant operation; the circulant structure emerges itself from this requirement. Finally, we will state a textbook result about commutative matrices, which motivates our next steps:

Theorem 18. *Two matrices A, B are jointly diagonalizable iff they commute, i.e., $AB = BA$.*

This theorem implies that circulant matrices are diagonalized by the eigenvectors of the shift operator. Our next step is to look closer into how such eigenvectors look like.

12.3 Eigenvalues and eigenvectors of S^\top

 12.3 Eigenvalues and eigenvectors of S^\top

We first find the eigenvalues of S^\top :

$$\begin{aligned} S^\top v = \lambda v &\iff v_{i+1 \bmod n} = \lambda v_i \\ (S^\top)^2 v = \lambda^2 v &\iff v_{i+2 \bmod n} = \lambda^2 v_i \\ &\vdots \\ (S^\top)^n v = \lambda^n v &\iff v_{i+n \bmod n} = \lambda^n v_i \end{aligned}$$

From the last equation, we get $v_{i+n \bmod n} = v_i = \lambda^n v_i$, and since $v \neq 0$, it necessarily follows $\lambda^n = 1$, i.e., the eigenvalues are the roots of unity

$$\lambda_k = e^{i\frac{2\pi}{n}k}.$$

The eigenvectors follow straightforwardly from the expression $v_k = \lambda^k v_0$, where from the requirement that $\|v\| = 1$ it follows that $v_0 = \frac{1}{\sqrt{n}}$ (recall that eigenvectors are defined up to scale). The eigenvectors can be written as the following orthonormal matrix ($V^*V = VV^* = I$),

$$V = \frac{1}{\sqrt{n}} \begin{pmatrix} 1 & 1 & \dots \\ 1 & e^{i\frac{2\pi}{n}} & \dots \\ \vdots & \vdots & \\ 1 & e^{i\frac{2\pi}{n}(n-1)} & \dots \end{pmatrix}$$

whose columns are given as integer powers of the vector $v = (1 e^{i\frac{2\pi}{n}} e^{i\frac{2\pi}{n}2} \dots e^{i\frac{2\pi}{n}(n-1)})^\top$. The matrix V gives rise to the *discrete Fourier transform (DFT)*, defined as

$$\hat{x}_k = (V^*x)_k = \frac{1}{\sqrt{n}} \sum_{\ell=0}^{n-1} x_\ell e^{-i\frac{2\pi}{n}k\ell}.$$

DFT is ubiquitously used in signal processing for reasons that will become clear in the following. Note that though V is a dense $n \times n$ matrix, its special structure allows performing multiplications of the form Vx (forward FFT) and $V^*\hat{x}$ (inverse DFT) in $\mathcal{O}(n \log n)$ complexity instead of $\mathcal{O}(n^2)$ using the *fast Fourier transform (FFT)* algorithm. The key idea of FFT (or a particular method known as the Cooley-Tukey or radix-2 algorithm) is to recursively split the vector into odd and even parts, reusing the product results.

12.4 Eigenvalues of circulant matrices

Since $C(x)$ and S^\top commute, we conclude that $C(x)$ is diagonalized by the DFT, i.e.

$$C(x) = V \begin{pmatrix} \mu_0 & & \\ & \ddots & \\ & & \mu_{n-1} \end{pmatrix} V^*,$$

12.4 Eigenvalues of circulant matrices

where $\{\mu_0, \dots, \mu_{n-1}\}$ are the eigenvalues of $C(x)$, which we will now express. For the k th eigenvector and eigenvalues we obtain

$$\begin{pmatrix} x_0 & x_{n-1} & & x_1 \\ x_1 & x_0 & & x_2 \\ \vdots & & \ddots & \vdots \\ x_{n-1} & x_{n-2} & \dots & x_0 \end{pmatrix} \begin{pmatrix} 1 \\ e^{i\frac{2\pi}{n}k} \\ \vdots \\ e^{i\frac{2\pi}{n}k(n-1)} \end{pmatrix} = \mu_k \begin{pmatrix} 1 \\ e^{i\frac{2\pi}{n}k} \\ \vdots \\ e^{i\frac{2\pi}{n}k(n-1)} \end{pmatrix},$$

where we omitted the $\frac{1}{\sqrt{n}}$ normalization on both sides. The first row yields the desired expression:

$$\begin{aligned} \mu_k &= x_0 + x_{n-1}e^{i\frac{2\pi}{n}k} + \dots + x_1e^{i\frac{2\pi}{n}k(n-1)} \\ &= x_0 + x_1e^{-i\frac{2\pi}{n}k} + \dots + x_{n-1}e^{-i\frac{2\pi}{n}k(n-1)} \\ &= \sum_{\ell=0}^{n-1} x_\ell e^{-i\frac{2\pi}{n}k\ell} = \hat{x}_k. \end{aligned}$$

We conclude that circulant matrix can be written in the form

$$C(x) = V \begin{pmatrix} \hat{x}_0 & & & \\ & \hat{x}_1 & & \\ & & \ddots & \\ & & & \hat{x}_{n-1} \end{pmatrix} V^*,$$

which gives a simple and efficient recipe for performing convolutions $x \circledast y = C(x)y$ in the Fourier domain:

1. Compute FFT: $\hat{y} = V^*y$ with $\mathcal{O}(n \log n)$ complexity.
2. Apply filter by means of element-wise product: $\hat{x} \circ \hat{y}$ (equivalent to multiplication by a diagonal matrix) in $\mathcal{O}(n)$ complexity.
3. Compute inverse FFT: $x \circledast y = V(\hat{x} \circ \hat{y})$ in $\mathcal{O}(n \log n)$ complexity.

Remark. Note that the discrete Laplacian operator

$$\Delta = \begin{pmatrix} -2 & 1 & & 1 \\ 1 & -2 & 1 & \\ \vdots & & \ddots & \vdots \\ 1 & & & 1 & -2 \end{pmatrix}$$

is also a circulant matrix, hence diagonalized by the DFT. We can therefore define convolutions as linear operations commuting with the Laplacian.