# The 2D Heat Equation
## — Analysis and Numerical Approaches

Jiajun Wang and Jiongyi Wang

Last Updated on: April 4, 2022

# Contents

# 1 Formulation and Primary Exploration of Problem

## 1.1 Physical Interpretation

We want to study how heat is conducted in a metal rod of length $L$ over time.One end of the metal rod is at and the other end is at. The length of the metal rod is much greater than its cross-sectional radius, so we can think of heat conduction as a function of $x$ and $t$.

Assuming the specific heat capacity of the metal rod is known, if we can find a function of temperature, we can know how the heat diffuses.

The rod is assumed to be adiabatic along its length, so it can only absorb or dissipate heat through the ends. This means that the temperature distribution depends only on the following three factors:

1. Initial temperature distribution, $T(x, 0)$this is called initial condition.
2. The temperature at both ends of the metal rod $T(0, t), T(L, t)$is called boundary conditions.
3. The law of heat transfer from one point to another in the metal rod. The heat equation is a mathematical representation of this physical law.

We assume that the initial boundary value of the solved heat equation $T(0, t) = T(L, t) = 0$

The problem of solving partial differential equations for a specific set of initial and boundary conditions is called an initial boundary value problem.

The heat equation can be derived from the conservation of energy: the time rate of change of the heat stored at a point on the metal rod is equal to the net heat flux into that point. This process fits the continuity equation. If $Q$ is the heat at various points and $V$ is the vector field of heat flow, then:

$$\frac{\partial Q}{\partial t} + \nabla V = 0$$

According to the second law of thermodynamics, if two identical bodies are in thermal contact, one being hotter than the other, then heat must flow from the hotter body to the cooler body at a rate proportional to the temperature difference. Therefore, V is proportional to the negative gradient of temperature, so: $V = -kT$ where k is the thermal conductivity of the metal. In one dimension, $Q = \rho c T$, $k$, $\rho$, and $c$ are the thermal conductivity, density, and specific heat capacity of the metal, respectively. Substituting into the expressions for $V$ and $Q$ yields the heat equation:

$$\frac{\partial T}{\partial t} - \frac{k}{\rho c}\frac{\partial^2 T}{\partial x^2} = 0$$

Now we give a formal formulation for the heat equation:

Let $\Omega \subset \mathbb{R}^d$ be an open set with boundry $\Gamma := \partial\Omega$, set $\Omega_T = \Omega \times ]0, T[$, $\Gamma_T := \Gamma \times ]0, T[$, $\Gamma_T$ is called the *lateral* boundary of the cylinder $\Omega_T$.

Consider the heat equation with $L$-periodic boundary condition:

$$\begin{cases} \partial_t u - k\Delta u &= f(x, t) & \text{in } \Omega_T \\ u(x, t) &= u(x + L, t) & \text{on } \Gamma_T \\ u(x, 0) &= u_0(x) & \text{on } \Omega \times \{t = 0\} \end{cases} \tag{1.1}$$

N.B. where $k > 0$ is a 'diffusion coefficient'. However, since the constant can be scaled out by defining a rescaled time $\tau = t/k$ to get

$$u_\tau - \Delta u = f$$

Hence we could simplify the formulation as

$$\begin{cases} \partial_t u - \Delta u & = f(x,t) & \text{in } \Omega_T \\ u(x,t) & = u(x+L,t) & \text{on } \Gamma_T \\ u(x,0) & = u_0(x) & \text{on } \Omega \times \{t = 0\} \end{cases} \tag{1.2}$$

In high dimensional case, $\Delta u = \Delta_x u = \sum_1^d \frac{\partial^2 u}{\partial x_i^2} = \text{div}\,(\text{grad } u)$. Finally, in PDE context, we often note $\frac{\partial u}{\partial t}$ as $u_t$ and $u_{xx} := \frac{\partial^2 u}{\partial x^2} = \Delta u$ if no ambiguity occurs.
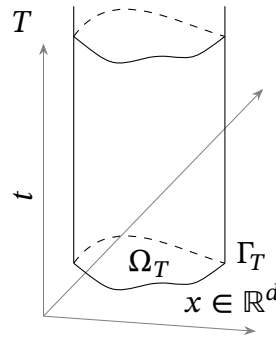


Figure 1: Region $\Omega_T$

## The Modeling Context: An Excursion to IC and BCs

1. **Dirichlet**: The temperature is fixed at the boundary
2. **Neumann**: The end is insulated (no heat enters or escapes).
3. **Robin**: Some heat enters or escapes, with an amount proportional to the temperature:

$$\alpha u = -\beta \frac{\partial u}{\partial n}$$

where $\frac{\partial u}{\partial n} := \nabla u \cdot \vec{n}$. For the area $\Omega$ whether heat enters or escapes the system depends on the boundary $\Gamma = \partial \Omega$. The heat flux $-\beta \frac{\partial u}{\partial n}$ is to the right if it is positive, so at the point $a$ belongs to boundary, heat enters the system when $\alpha > 0$ and leaves when $\alpha < 0$.

4. **Periodic**: The temperature is periodic at the boundary.

The same interpretations apply when the equation is describing diffusion of some other quantity. Non-homogeneous boundary conditions can be imposed, for instance

$$u(a,t) = g(a,t), \quad a \in \Gamma_T$$

which might be used to model the ambient temperature increasing with time.

> **Remark 1** (Robin's BC): All BCs could be expressed as Robin's condition
>
> $$\alpha u(a) + \beta \frac{\partial u}{\partial n}(a) = g, \quad a \in \partial\Omega$$
>
> Dirichlet's condition means that $\beta = 0$ and Neumann's condition means that $\alpha = 0$.
>
> **Remark 2:** At each point of $\Gamma$, the BCs could be different.

## 1.2 Fundamental Solution for Heat Equation

We observe that the heat equation involves one derivative with respect to the time variable $t$ ,and two derivatives with respect to the space variable $x_i (i = 1, .....n)$. At the same time ,we can observe that if $u$ is a solution of the heat equation , $u(\lambda x, \lambda^2 t)$ for $\lambda \in \mathbb{R}$ is also a solution. So the $\frac{r^2}{t}, (r = |x|)$ is very important for the heat equation. We need to find a solution as $u(x, t) = v(\frac{|x|^2}{t})(t > 0, x \in \mathbb{R})$.

But here we want an idea for which we can more easily find a solution [cf. Eva10, p.45] $u(x, t) = \frac{1}{t^a} v(\frac{x}{t^b})$ and we have $x$ in $\mathbb{R}, t > 0$.

For this equation, the result 1 is easily obtained by a simple calculation.

$$\lambda^a u(\lambda^b, \lambda t) = \frac{\lambda^a}{(\lambda t)^a} v(\frac{\lambda^b x}{(\lambda t)^b}) = \frac{1}{t^a} v(\frac{x}{t^b}) = u(x, t)(result 1)$$

Let $\lambda = \frac{1}{t}$ and $y = \lambda^b x$ , we can get $v(y) = u(y, 1)$. We substitute $u(x, t)$ into the first expression of the heat equation and get

$$at^{-a-1} v(y) + bt^{-a-1} y.Dv(y) + t^{-a-2b} \Delta v(y) = 0$$

To simplify the equation, we assume $v(y) = w(|y|)$, and $w : \mathbb{R} \to \mathbb{R}$. Then the above formula becomes:

$$at^{-a-1} v(y) + bt^{-a-1} y.Dv(y) + t^{-a-2b} \Delta v(y) = aw + \frac{1}{2} rw' + w'' + \frac{n-1}{r} w' = 0$$

if we let a= $\frac{n}{2}$,we will get

$$\frac{n}{2} w + \frac{1}{2} rw' + w'' + \frac{n-1}{r} w' = (r^{n-1} w' + \frac{1}{2} r^n w)' = 0$$

$\Rightarrow r^{n-1} w' + \frac{1}{2} r^n w = C$ and we assume that $\lim_{r \to 0} w = 0, \lim_{r \to 0} w' = 0$, and we can derive $C = 0$ and $w' = -\frac{1}{2} rw$. Thus, we have $w = de^{-\frac{r^2}{4t}}$ .

Finally, with $u(x, t) = \frac{1}{t^{\frac{n}{2}}} v(\frac{x}{t^{\frac{1}{2}}})$, we deduce $u(x, t) = \frac{d}{t^{\frac{n}{2}}} e^{-\frac{|x|^2}{4t}}$

So we define the basic solution of the heat equation as

✠ **Definition 1.1:**  *The function*

$$\Phi(x,t) = \begin{cases} \dfrac{1}{(4\pi t)^{d/2}} e^{-\frac{|x|^2}{4t}} & x \in \mathbb{R}^d, t > 0 \\ 0 & x \in \mathbb{R}^d, t < 0 \end{cases} \tag{1.3}$$

*is called the fundamental solution of the heat equation.*

The choice of the normalizing constant $(4\pi)^{-\frac{n}{2}}$ is dictated by

$$\int_{\mathbb{R}^n} \Phi(x,t)dx = \frac{1}{(4\pi t)^{\frac{n}{2}}} \int_{\mathbb{R}^n} e^{-\frac{|x|^2}{4t}} dx$$

$$= \frac{1}{\int_{\mathbb{R}^n}} e^{-|z|^2} dz = \frac{1}{\pi^{\frac{n}{2}}} \prod_n^{i=1} \int_{-\infty}^{+\infty} e^{-z_i^2} dz_i$$

$$= 1$$

Next, we need to verify the correctness of this result. We need to calculate $\frac{\partial \Phi(x,t)}{\partial t}$ and $\frac{\partial^2 \Phi(x,t)}{\partial x_i^2}$ separately, and then substitute the results into the heat equation for verification. So, we can get the following result:

$$\frac{\partial \Phi(x,t)}{\partial t} = \frac{1}{(4\pi)^{\frac{n}{2}}}\left((-\frac{n}{2})(t^{-\frac{n}{2}-1})e^{-\frac{|x|^2}{4t}} + \frac{1}{t^{\frac{n}{2}}}e^{-\frac{|x|^2}{4t}}\frac{|x|^2}{4t}\right)$$

$$= \frac{e^{-\frac{|x|^2}{4t}}}{(4\pi)^{\frac{n}{2}}}\left((-\frac{n}{2})\frac{1}{t^{\frac{n}{2}+1}} + \frac{|x|^2}{4t^{\frac{n}{2}+2}}\right) \tag{1.4}$$

$$\frac{\partial \Phi(x,t)}{\partial x_i} = \frac{e^{-\frac{|x|^2}{4t}}}{(4\pi t)^{\frac{n}{2}}}\left(-\frac{1}{2t}x_i\right)$$

$$\frac{\partial^2 \Phi(x,t)}{\partial x_i^2} = \frac{e^{-\frac{|x|^2}{4t}}}{(4\pi t)^{\frac{n}{2}}}\left(\frac{1}{4t^2}x_i^2\right) - \frac{1}{2t}\frac{e^{-\frac{|x|^2}{4t}}}{(4\pi t)^{\frac{n}{2}}}$$

$$\sum_{i=0}^{n} \frac{\partial^2 \Phi(x,t)}{\partial x_i^2} = \frac{e^{-\frac{|x|^2}{4t}}}{(4\pi t)^{\frac{n}{2}}}\left(\frac{1}{4t^2}|x|^2\right) - \frac{n}{2t}\frac{e^{-\frac{|x|^2}{4t}}}{(4\pi t)^{\frac{n}{2}}}$$

$$= \frac{e^{-\frac{|x|^2}{4t}}}{(4\pi)^{\frac{n}{2}}}\left((-\frac{n}{2})\frac{1}{t^{\frac{n}{2}+1}} + \frac{|x|^2}{4t^{\frac{n}{2}+2}}\right) \tag{1.5}$$

$$\tag{1.6}$$

By substituting (2) and (3) into the heat equation, we can find that

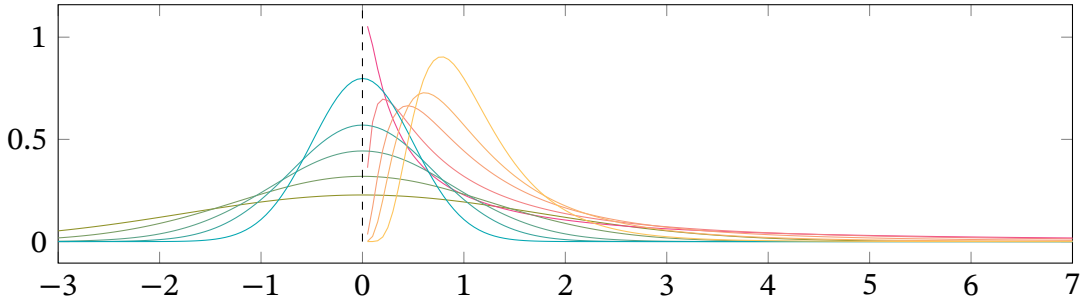$$\frac{\partial u}{\partial t}(x,t) + \frac{\partial^2 u}{\partial x^2}(x,t) = 0$$

Figure 2: example for drawing fundamental solution with varing $t$ [not finished]

## 1.3 Spectrum Theory and Spectral Analysis

In this part, for analysis preliminaries, [Fol99, ch.05-06] offers a panorama on the theory of $L^p$ spaces.

Given $E, F$ two normed vector spaces, an operator $T \in \mathcal{L}(E; F)$ is said to be copmact if the image of unit ball in $E$ under $T$, i.e. $T(B_E)$ is relatively compact in $F$. We call an operator $T$ is of finite rank, if $\dim \operatorname{Im} T < \infty$.

### 1.3.1 Riesz-Fredholm theory

▶ **Lemma 1.2** (Riesz)**:** *Let $E$ be a normed vector space (not necessary complete), $M \subsetneq E$ a proper closed linear subspace, then $\forall \varepsilon > 0$, $\exists u \in E$ s.t. $\|u\| = 1$ and $d(u, M) \geq 1 - \varepsilon$.*

▶ **Theorem 1.3:** *Let $E$ be a normed vector space with compact unit ball $B_E$, then $E$ is finite-dimensional.*

▶ **Theorem 1.4** (Fredholm alternative)**:** *Let $T \in \mathcal{L}(E)$ be a compact operator, then*

- $\operatorname{Ker}(I - T)$ *is finite-dimensional.*
- $\operatorname{Im}(I - T)$ *is closed. More precisely,* $\operatorname{Im}(I - T) = \operatorname{Ker}(I - T')^{\perp}$
- $\operatorname{Ker}(I - T) = \{0\} \iff \operatorname{Im}(I - T) = E$
- $\dim \operatorname{Ker}(I - T) = \dim \operatorname{Ker}(I - T')$

▷ *Proof:* Admitted. [cf. Bre11, p.160-162] □

✠ **Definition 1.5** (resolvent set, spectrum and eigenvalue)**:** *Let $T \in \mathcal{L}(E)$, the resolvent set, denoted by $\rho(T)$, is defined by*

$$\rho(T) := \{\lambda \in \mathbb{C}; (T - \lambda I) : E \to E \text{ is bijective}\}$$

*The spectrum, denoted by $\operatorname{Spec}(T)$, is the complement of the resolvent set, i.e., $\operatorname{Spec}(T) = \mathbb{C} \backslash \rho(T)$. A complex number $\lambda$ is said to be an eigenvalue of $T$ if $\operatorname{Ker}(T - \lambda I) \neq \{0\}$. The set of eigenvalues of $T$ is denoted by $EV(T)$. The space $\operatorname{Ker}(T - \lambda I)$ is called the eigenspace of $T$, the elemet in it is called eigenvector.*

**Remark 3:** If $E$ is a Banach space, the open mapping theorem tells us, the bijectivity of $T$ equals that $T^{-1} \in \mathcal{L}^{-1}(E)$. Actually we have following consequnce :

▶ **Proposition 1.6:**

- *If $T \in \mathcal{L}(E)$ and $\|I\text{-}T\| < 1$ where $I$ is the identity operator, then $T$ is invertible, the series $\lim_{n\to\infty} \sum_{n=0}^{\infty}(I-T)^n = T^{-1}$ in $\mathcal{L}(E)$.*
- *The set of invertible operators in $\mathcal{L}(E)$, denoted as $GL(E)$, is an open set in $\mathcal{L}(E)$, and $GL(E) \to GL(E); T \mapsto T^{-1}$ is continuous. More precisely, if $S \in GL(E)$ and $\|T\text{-}S\| < \left\|T^{\text{-}1}\right\|^{-1}$, then $S \in GL(E)$.*

▶ **Theorem 1.7** (Gelfand)**:**

- *We have $\|T^n\|^{1/n} \xrightarrow{n\to\infty} \inf_n \|T^n\|^{1/n}$, we call this limite, denoted by $r(T)$, the spectral radius of $T$. Moreover, $r(T) \leq \|T\|$, and $\forall \lambda \in \mathrm{Spec}(T), |\lambda| \leq r(T)$. In particular, $\mathrm{Spec}(T)$ is a compact set in $\mathbb{C}$.*
- *For all $T \in \mathcal{L}(E)$, we have $\mathrm{Spec}(T) \neq \varnothing$. Moreover*

$$r(T) = \max_{\lambda \in \mathrm{Spec}(T)} \{|\lambda|\}$$

▷ *Proof:*  [cf. Lax02, p195-197]  □

✠ **Definition 1.8** (adjoint, self-adjoint)**:**  *Let $A : \mathrm{Dom}\,(A) \subset E \to F$ be an unbounded linear operator that is densely defined. We shall introduce an unbounded operator $A' : \mathrm{Dom}\,(A') \subset F' \to E'$ as follows:*

$$\mathrm{Dom}\,(A') := \{v \in F' : \exists c \geq 0 \text{ s.t. } |\langle v, Au \rangle| \geq c\,\|u\|, \quad \forall u \in \mathrm{Dom}\,(A)\}$$

$$_{F'}\langle v, Au \rangle_F = {}_{E'}\langle A'v, u \rangle_E, \quad \forall u \in \mathrm{Dom}\,(A), \forall v \in \mathrm{Dom}\,(A')$$

*A bounded operator $T$ is said to be self-adjoint if $T' = T$.*

▶ **Theorem 1.9:**  *Suppose that $H$ is a separable Hilbert space, $T$ is a compact self-adjoint operator, then there exists a Hilbert basis composed of eigenvectors of $T$.*

Our last statement is a fundamental result. It asserts that every compact self-adjoint operator may be diagonalized in some suitable basis.

## 1.3.2  Eigenfunctions and spectral decomposition

Now, we have sufficient tools to proceed the spectral analysis of heat equation. (More generally, the spectral analysis could be applied to other types of PDE [cf. Bre11, ch.08-09; Lax02, ch.33-36])

✠ **Definition 1.10** (Sobolev spaces, distribution)**:**

$$W^{m,p}(\Omega) := \left\{u \in L^p(\Omega) : \partial^\alpha u \in L^p(\Omega), \forall \alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{R}_+^d \text{ and } 1 \leq |\alpha| \leq m\right\}$$

*For index $\alpha \in \mathbb{R}_+^d$, we note $|\alpha| = \sum_1^d \alpha_i$. The norm $\|\cdot\|_{W^{m,p}(\Omega)}$ defined by*

$$\|u\|_{W^{m,p}(\Omega)} = \left(\sum_{|\alpha|=0}^m \|\partial^\alpha u\|_{L^p(\Omega)}^p\right)^{\frac{1}{p}}$$

*makes the Sobolev space $W^{m,p}(\Omega)$ complete.*

*We simply note $H^m(\Omega) := W^{m,2}(\Omega)$, since it's a Hilbert space.*

*At last, we define $H_0^m(\Omega) := \overline{\mathcal{D}(\Omega)}^{H^m(\Omega)}$, that means the closure of $\mathcal{D}(\Omega)$ in $H^m(\Omega)$.*

*Where $\mathcal{D}(\Omega) = C_c^\infty(\Omega)$ called the set of test functions. Its dual space $\mathcal{D}'(\Omega)$ is called distribution.*

*$H_0^1(\Omega)$ need not to inherit the norm from $H^1(\Omega)$, there is an equivalent norm inducted by the inner product:*

$$\langle v, u \rangle_{H_0^1(\Omega)} = \langle \nabla u, \nabla v \rangle_{L^2(\Omega)}$$

The notion of distriburion generalized the notion of function. We could find that the Dirac mass at $x = a$: $\delta(x - a)$ is not a function, however, it's a distriburion. Important example: $L_{loc}^1(\Omega)$ is a distribution. For more details, [cf. Gos20, ch.03-05].

▶ **Theorem 1.11** (the spectrum of Laplacian operator)**:** *Suppose $T : L^2(\Omega) \to L^2(\Omega); f \mapsto u$, where $u$ is the weak solution (i.e. solution in $H_0^1(\Omega)$) of*

$$\begin{cases} -\Delta u_f = f & \text{in } \Omega \\ u_f = 0 & \text{on } \partial\Omega \end{cases} \tag{1.7}$$

*Then $T$ is compact and self-adjoint. Moreover $T$ is positive defined.*

▷ *Proof:* $u_f$ is characterized by $u_f \in H_0^1(\Omega), \forall \varphi \in H_0^1(\Omega), \int_\Omega \nabla u_f \nabla \varphi = \int_\Omega f\varphi$



$T = i \circ S$, where $S$ is linear continuous and $i : H_0^1(\Omega) \to L^2(\Omega)$ is compact injection (it's a result of Reillich-Kondrachov's Theorem, [cf. Bre11, ch.9.3, p.285], we don't discuss the interpolation of Sobolev space here)

To prove that $T$ self-adjoint and positive is relatively easy, it's direct consequnce of properties of $L^2(\Omega)$. □

Thus, by theorem 1.9, there exists a Hilbert basis in $L^2(\Omega)$ consists of the eigenvalues of $T$: assume that $\Omega \subset \mathbb{R}^d$ is a bounded open set, then there exist a Hilbert basis $\{e_n\}_n$ of $L^2(\Omega)$ s.t. $e_n \in H_0^1(\Omega) \cap C^\infty(\Omega), \forall n$ and a sequence $\{\lambda_n\}_n$ of real numbers with $\lambda_n > 0, \forall n$ and $\lambda_n \to +\infty$ s.t.

$$-\Delta e_n = \lambda_n e_n \quad \text{in } \Omega$$

We say that $\{\lambda_n\}_n$ are the eigenvalues of $-\Delta$ ( with Dirichlet boundary condition) and the $\{e_n\}_n$ are the associated eigenfunctions.

▷ *Proof:* □

▶ **Corollary 1.12:** $\left\{ \frac{1}{\sqrt{\lambda_n}} e_n \right\}_n$ *is a Hilbert basis for $H_0^1(\Omega)$ equipped with the inner product*

$$\langle v, w \rangle_{H_0^1(\Omega)} = \int_\Omega \nabla v \nabla w$$

.

▷ *Proof:*

$$\left\langle \nabla \frac{1}{\sqrt{\lambda_n}} e_n, \nabla \frac{1}{\sqrt{\lambda_m}} e_m \right\rangle_{L^2(\Omega)} = \frac{1}{\sqrt{\lambda_n}} \left\langle \sqrt{\lambda_n} \nabla T e_n, \nabla \frac{1}{\sqrt{\lambda_m}} e_m \right\rangle_{L^2(\Omega)}$$

$$= \frac{1}{\sqrt{\lambda_n}} \sqrt{\lambda_n \lambda_m} \int_\Omega e_n e_m$$

$$= \sqrt{\frac{\lambda_m}{\lambda_n}} \langle e_n, e_m \rangle_{L^2(\Omega)} = \sqrt{\frac{\lambda_m}{\lambda_n}} \delta_n^m$$

Then $\left\langle \nabla \sqrt{\lambda_n} e_n, \nabla \sqrt{\lambda_m} e_m \right\rangle = \sqrt{\frac{\lambda_m}{\lambda_n}} \delta_n^m \implies \left\{ \sqrt{\lambda_n} e_n \right\}_n$ is an orthonomal family in $H_0^1(\Omega)$.

Lastly, we should prove $\left\{ \sqrt{\lambda_n} e_n \right\}_n$ is a maximal family, i.e. it spans a dense party of $H_0^1(\Omega)$.

Let $v \in H_0^1(\Omega)$ s.t. $\forall n, \left\langle \sqrt{\lambda_n} \nabla e_n, \nabla v \right\rangle_{L^2(\Omega)} = 0 \implies \frac{1}{\sqrt{\lambda_n}} \langle \nabla T e_n, \nabla v \rangle_{L^2(\Omega)} = 0 \implies \langle e_n, v \rangle_{L^2(\Omega)} = 0$.

Since $\{e_n\}_n$ is a Hilbert basis of $L^2(\Omega)$, then $v = 0$ in $L^2(\Omega)$, naturally $v = 0$ in $H_0^1(\Omega)$.

We have hereby proved that $\left( \overline{\text{span}(\{e_n\})}^{H_0^1(\Omega)} \right)^\perp = \{0\}$ in $H_0^1(\Omega)$, then $\left\{ \sqrt{\lambda_n} e_n \right\}_n$ is a maximal family.

$\square$

However, our problem is more complicated than the case above, that's because the restriction of initial condition (IC) and boundary conditions (BCs) may not allow the solution belong to $H_0^m(\Omega)$, which implies only a zero Dirichlet BC. Thus, the choice of a proper working space is the first and the crucial matter to think about.

---

### Separation of Variables and Superpostition Principle: A History

When Joseph Fourier contemplate the heat equation
As we did before, sometimes we note $u_{xx}$ as $\Delta u$, because the Laplacian operator is a linear operator on $u$ with respect to $x$ — actually, the heat equat also belongs to linear PDE, which refers that

$$u_t = -L[u] + f(x, t) \tag{1.8}$$

A linear, homogeneous PDE obeys the superposition principle: $u_1, u_2$ are solutions $\implies c_1 u_1 + c_2 u_2$ is a solution for all scalars $c_1, c_2 \in \mathbb{R}$. The concepts of linearity and homogeneity also apply to boundary conditions, in which case the variables are evaluated at specific points.
For example, here lists the linear operators for some basic linear PDEs, i.e. heat equation, wave equation, Poisson's equation.

$$L_h = \partial_t - \partial_{xx}, \quad L_w = \partial_{tt} - \partial_{xx}, \quad L_P = \nabla^2$$

an example of non-linear PDE:

$$u_t + u u_x = u_{xx}$$

However, $L_h$ is generally not a compact and self-adjoint operator, which means we should consider rather the Laplacian operator $\Delta$ and the basis is a series of functions involves only $x$.

To have an eigenfunction of the operator $L$, we must prescribe $\Omega$ and associated boundary conditions. Generally, we could sketch out the proper working space $V$ as

$$V = \{f \in H^m(\Omega) : f \text{ satisfies BCs}\} \tag{1.9}$$

The eigenfunctions we need is to solve the eigenvalue problem

$$L[\phi] = \lambda\phi, \quad \phi \in V$$

By theorem 1.9, there is a sequence of eigenfunctions $\{\phi_n\}$ with eigenvalues $\{\lambda_n\}$ that form an orthogonal basis for the space $V$.

Now at each fixed time $t$, the function $u(x,t)$, regarded as a function of $x$, lies inside the space $V$. It follows that there are coefficients $a_n(t)$ such that

$$u(x,t) = \sum_{n=0}^{\infty} a_n(t)\phi_n(x) \tag{1.10}$$

Our objective now is to determine the functions $a_n(t)$. We write the heat source $f$ in terms of eigenfunctions:

$$f(x,t) = \sum_{n=0}^{\infty} f_n(t)\phi_n(x) \tag{1.11}$$

Now substitute this 1.11 and the eigenfunction expansion for $u$ (equation 1.11) into the PDE (30) to obtain

$$\underbrace{\sum_{n=0}^{\infty} a_n'(t)\phi_n(x)}_{u_t} = \underbrace{-\sum_{n=0}^{\infty} a_n(t)\lambda_n\phi_n(x)}_{-L[u]} + \underbrace{\sum_{n=0}^{\infty} f_n(t)\phi_n(x)}_{f} \tag{1.12}$$

In detail, the second term was found using the eigenfunction property and linearity of $L$:

$$\begin{aligned}
L[u] &= L\left[\sum_{n=0}^{\infty} a_n(t)\phi_n(x)\right] \\
&= \sum_{n=0}^{\infty} a_n(t)L[\phi_n(x)] \\
&= \sum_{n=0}^{\infty} a_n(t)\lambda_n\phi_n(x)
\end{aligned}$$

Note that since $a_n(t)$ is only a function of $t$, it is constant as far as $L$ is concerned so by linearity, $L[a_n(t)\phi_n(x)] = a_n(t)L[\phi_n(x)]$, where we have used the fact that $L$ is linear to move.

Now we gather equation 1.12 together under the basis $\{\phi_n\}$:

$$0 = \sum_{n=0}^{\infty} [\underbrace{a_n'(t) + \lambda_n a_n(t) - f_n(t)}_{\Psi_n(t)}]\phi_n(x) = \sum_{n=0}^{\infty} \Psi_n(t)\phi_n(x)$$

Since the $\{\phi_n(x)\}$'s are a basis, the coefficient of each basis function must be zero at all times $t$ (otherwise, the $\{\phi_n(x)\}$'s would be linearly dependent at some $t$), i.e. $\Psi_n(t) \equiv 0$. It follows that for each $n$,

$$a_n'(t) + \lambda_n a_n(t) = f_n(t), \quad \forall n > 0, t \in ]0, T] \tag{1.13}$$

This equation is a first-order linear ODE for $a_n(t)$ that is easy to solve.

To complete our problem, the last missing piece is the initial condition, i.e. $a_n(0)$.

Recall we set $u(x, 0) = u_0(x)$. Write the IC in terms of eigenfunctions:

$$u_0(x) = \sum_{n=0}^{\infty} u_{0,n} \phi_n(x)$$

For the solution 1.10 satisfies the IC, we need the constants $\gamma_n$

$$\underbrace{\sum_{n=0}^{\infty} a_n(0) \phi_n(x)}_{u(x,0)} = \underbrace{\sum_{n=0}^{\infty} u_{0,n} \phi_n(x)}_{u_0(x)}$$

Again, since the $\{\phi_n(x)\}$'s are absis, the two sums must be equal term-by-term, so

$$a_n(0) = u_{0,n}, \quad \forall n$$

Finally, this condition and equation 1.13 lets us solve for a unique $a_n(t)$ (as we get a first order IC problem), which completes the process.

## Sturm–Liouville Eigenvalue Problems

We have found the method of separation of variables to be quite successful in solving some homogeneous partial differential equations with homogeneous boundary conditions.

In all examples we have analyzed so far, the boundary value problem that determines the needed eigenvalues (separation constants) has involved the simple ordinary differential equation

$$\phi''(x) + \lambda \phi(x) = 0 \tag{1.14}$$

$$\frac{d}{dx}\left(p \frac{d\phi}{dx}\right) + q\phi + \lambda \sigma \phi = 0 \tag{1.15}$$

Explicit solutions of this equation determined the eigenvalues $\lambda$ from the homogeneous boundary conditions.

▶ **Proposition 1.13:** *For any regular Sturm-Liouville problem, all of the following theorems are valid:*

1. *All the eigenvalues $\lambda$ are real.*
2. *There exist an infinite number of distinct eigenvalues:*

$$\lambda_1 < \lambda_2 < \cdots < \lambda_n < \cdots \to +\infty$$

3. *Corresponding to each eigenvalue $\lambda_n$, there is an eigenfunction, denoted $\phi_n(x)$ (which is unique to within an arbitrary multiplicative constant)*

11

4. The eigenfunctions $\phi_n(x)$ form a "complete" set, meaning that any piecewise smooth function $f(x)$ can be represented by a generalized Fourier series of the eigenfunctions:

$$f(x) \sim \sum_{n=0}^{\infty} a_n \phi_n(x)$$

5. Eigenfunctions belonging to different eigenvalues are orthogonal relative to the weight function $\sigma(x)$. In other words,

$$\int_{\Omega} \phi_m(x)\phi_n(x)\sigma(x)dx = 0, \quad \lambda_m \neq \lambda_n$$

6. Any eigenvalue can be related to its eigenfunction by the Rayleigh quotient

$$\lambda = \frac{-\int_{\Gamma} p\phi\frac{\partial\phi}{\partial n}d\Gamma + \int_{\Omega}[p(\frac{d\phi}{dx})^2 - q\phi^2)]}{\int_{\Omega} \phi^2\sigma}$$

## 1.4 Maximum Principle

▶ **Theorem 1.14** (Strong maximum principle)**:** *Assume that $u \in C^2(\Omega) \cap C^1(\overline{\Omega})$ solve the heat equation in $\Omega_T$, then*

$$\max_{\overline{\Omega}_T} u = \max_{\partial\Omega_T} u$$

▶ **Theorem 1.15:** *Assume that $u_0 \in L^2(\Omega)$, $u$ is the solution of the equation 1.2, then we have, for all $(x, t) \in \Omega_T$*

$$\min\left\{0, \inf_{\Omega} u_0\right\} \leq u(x, t) \leq \max\left\{0, \sup_{\Omega} u_0\right\}$$

▷ *Proof:* Instead of a classical mean value formula [cf. Eva10, ch.02.3.2-02.3.3], we use Stampacchia's truncation method. Set

$$K = \max\left\{0, \sup_{\Omega} u_0\right\}$$

and assume that $K < +\infty$. Fix a function $G$ s.t.

- $|G'(s)| \leq M, \quad \forall s \in \mathbb{R}$
- $G$ is strictly increasing on $]0, \infty[$
- $G(s) = 0, \quad \forall s \leq 0$

and let

$$H(s) = \int_{0}^{s} G(r)\,dr, \quad s \in \mathbb{R}$$

Easy to check the function $\varphi$ defined by

$$\varphi(t) = \int_\Omega H(u(x,t) - K)\,\mathrm{d}x$$

has the following properties:

- $\varphi \in C([0,\infty[;\mathbb{R}), \varphi(0) = 0, \varphi \geq 0$ on $[0,\infty[$
- $\varphi \in C^1(]0,\infty[;\mathbb{R})$
- 

$$
\begin{aligned}
\varphi'(t) &= \int_\Omega G(u(x,t) - K)\partial_t u(x,t)\,\mathrm{d}x \\
&= \int_\Omega G(u(x,t) - K)\Delta u(x,t)\,\mathrm{d}x \\
&= -\int_\Omega G'(u - K)|\nabla u|^2\,\mathrm{d}x \\
&\leq 0
\end{aligned}
$$

since $G(u(x,t) - K) \in H^1_0(\Omega)$ for every $t > 0$, it follows that $\varphi \equiv 0$ and thus, $\forall t > 0, u(x,t) \leq K$ a.e. on $\Omega$. $\qquad\square$

▶ **Theorem 1.16:** *Assume $u \in C(\overline{\Omega} \times [0,T])$, $u$ is of class $C^1$ in $t$ and of class $C^2$ in $x \in \Omega \times ]0,T[$, $\partial_t u - \Delta_x u \leq 0 \in \Omega \times ]0,T[$, then*

$$\max_{\overline{\Omega}\times[0,T]} u = \max_{\mathfrak{P}} u$$

*where $\mathfrak{P} = (\overline{\Omega} \times \{0\}) \cup (\Gamma \times ]0,T[)$ is called the **parabolic boundary** of the cylinder $\Omega \times ]0,T[$*

# 2 Numerical Approaches

## 2.1 Discrete and Fast Fourier Transform (DFT & FFT)

In section 1.3.2 we have exhibited a complete procedure to solve heat equation in using eigenfunction's method, however, it's impossible to calculate a basis and corresponding coefficients presented by infinite series in computer. Generally we calculate a sufficient large number of finite items to approach the infinite series.

Let's start with a concrete example in 1D spacetime .

$$
\begin{cases}
\partial_t u - 2tu_{xx} &= e^{-t^2}\sin x \quad &\text{in } ]0,\pi[\times]0,1] \\
u(0,t) &= u(\pi,t) \quad &\text{on } \{0,\pi\}\times]0,1] \\
u(x,0) &= 3\sin 2x \quad &\text{on } ]0,\pi[\times\{t=0\}
\end{cases}
\tag{2.1}
$$

Eigenfunctions:

$$u_t = -2tL[u] + f(x,t)$$

where

$$L = \partial_{xx}, \quad f(x,t) = e^{-t^2} \sin x$$

Note that we cannot put $t$ in the definition of $L$, since $L$ must only involve derivatives in $x$. The eigenvalues/functions are

$$\lambda_n = n^2, \quad \phi_n(x) = \sin(nx), \quad n = 1, 2, \ldots$$

[cf. Sha03; Sch01, ch.08; Mal08, ch.03.3]
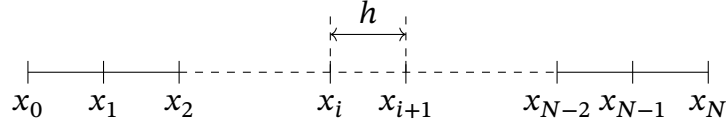because $\{e_k(t)\}_{k \in \mathbb{Z}} = \{e^{2i\pi kt}\}_{k \in \mathbb{Z}}$ is a Hilbert basis.



Figure 3: Equilong discretization in one dimension

The subsection 1.3 inspires us

$$\widehat{f}(k) = \int_0^\tau f(x)e^{-2i\pi kx}\, dx \xrightarrow{\text{discretization}} U_k = \frac{1}{N}\sum_{j=0}^{N-1} f(\frac{j}{N})e^{-2i\pi k\frac{j}{N}} \tag{2.2}$$

unequal lengths

### 2.1.1 FFT: interlacement

## 2.2 Finite Difference Method (FDM) for 1D Heat Equation

### 2.2.1 Explicit Euler's scheme

First, we determine the heat equation we need to consider

$$\begin{cases} \frac{\partial u}{\partial t}(x,t) + \frac{\partial^2 u}{\partial x^2}(x,t) &= f(x,t), \quad \text{in } \Omega \\ u(0,t) = u(L,t) &= g(t) \quad\quad \text{on } (0,T) \\ u(x,0) &= u_0(x) \quad \text{in } \Omega \end{cases} \tag{2.3}$$

The general idea of finite difference methods for evolution equations is that we replace the continuous variables in time and space with discrete points. In the case of an evolution equation, we need a spacetime grid. Let us thus be given two positive integers N and M. We set $h = \delta x = \frac{1}{N+1}$ and $X_n = nh$ for $n = 0, 1, 2, 3, \cdots, N+1$. We set $k = \delta t = \frac{T}{M+1}$ and $t_j = jk$ for $j = 0, 1, 2, 3, \cdots, M+1$. The parameter $h$ is called the space grid step and the parameter $k$ the time grid step, or time step. The grid points are the points $(x_n, t_j)$. we will let $N$ and $M$ go to infinity, or $h$ and $k$ go to 0.

The $u_n^j$ for the above values of $n$ and $j$, and it is hoped that $u_n^j$ will be an approximation of $u(x_n, t_j)$, that should become better and better as $N$ and $M$ are increased. The boundary condition can be enforced exactly by requiring that

$$u_0^j = u_{N+1}^j = 0$$

And we note that $j = 0, 1, 2, 3, \cdots, M + 1$. By the initial condition, we can acquire that

$$u_n^0 = u_0(x_n)$$

And we note that $n = 1, \cdots, N$.

In this way, we will find that all the points of the boundary are known.The only values that are left unknown at this stage are thus $u_n^j$ for $n = 1, 2, 3, \cdots, N$ and $j = 1, 2, 3, \cdots, M + 1$. We define that

$$U^j = \begin{pmatrix} u_1^j \\ u_2^j \\ \vdots \\ u_N^j \end{pmatrix} \in \mathbb{R}^N \tag{2.4}$$

Next, we need to calculate the remaining unknown points through the points of the known boundary. In the explicit Euler three point scheme, we can use the forward differential quotient approximation

$$\frac{\partial u}{\partial t}(x_n, t_j) \approx \frac{u(x_n, t_{j+1}) - u(x_n, t_j)}{k} \tag{2.5}$$

and for the second order space derivative,Using the same approximation, we can get:

$$\frac{\partial^2 u}{\partial x^2}(x_n, t_j) \approx \frac{\frac{u(x_{n+1}, t_j) - u(x_n, t_j)}{h} - \frac{u(x_n, t_j) - u(x_{n-1}, t_j)}{h}}{h} = \frac{u(x_{n+1}, t_j) - 2u(x_n, t_j) + u(x_{n-1}, t_j)}{h^2} \tag{2.6}$$

We substitute (3) and (4) into (1), we can get

$$\begin{cases} \frac{u_n^{j+1} - u_n^j}{k} - \frac{u_{n+1}^j - 2u_n^j + u_{n-1}^j}{h^2} = f_n^j & \text{for } n = 1, 2, 3, \cdots, N, j = 0, 1, 2, 3, \cdots, M, \\ u_n^0 = u_0(x_n) & \text{for } n = 1, \ldots, N \\ u_0^j = u_{N+1}^j = 0 & \text{for } j = 0, \cdots, M + 1 \end{cases} \tag{2.7}$$

So,we can rewrite the first N equations of the scheme in vector form as

$$\frac{U^{j+1} - U^j}{k} + A_h U^j = F^j \quad \text{for } j = 0, 1, 2, 3, \ldots, M. \tag{2.8}$$

where $A_h$ is the same $N \times N$ tridiagonal matrix

$$A_h = \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & -1 & 2 & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{pmatrix} \tag{2.9}$$

And we define that $F_j$ and the discrete initial condition are two vectors of $\mathbb{R}^N$.

$$F^j = \begin{pmatrix} f_1^j \\ f_2^j \\ \vdots \\ f_N^j \end{pmatrix} \in \mathbb{R}^N \qquad U_0 = \begin{pmatrix} u_0(x_1) \\ u_0(x_2) \\ \vdots \\ u_0(x_N) \end{pmatrix} \in \mathbb{R}^N \tag{2.10}$$

So,we can get the numerical scheme :

$$\begin{cases} U^{j+1} = (I - kA_h)U^j + kF^j \text{ for } j = 0, \cdots, M \\ U^0 = U_0 \end{cases}$$ (2.11)
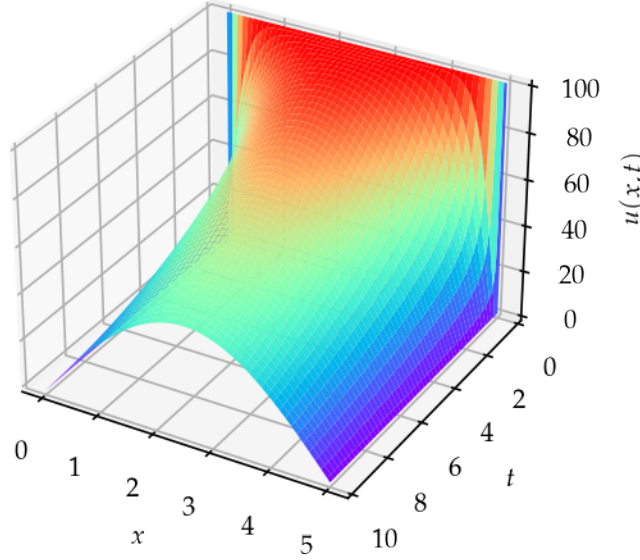


Figure 4: The numerical scheme for the 1D diffusion equation

### 2.2.2 Implicit Euler's scheme and Richardson method

For a full disciption of the scheme, [cf. Luc16, ch.08.2].
    The first example is the implicit or backward Euler three point scheme, which is associated with the backward differential quotient approximation of the time derivative

$$\frac{\partial u}{\partial t}(x_n, t_j) \approx \frac{u(x_n, t_j) - u(x_n, t_{j-1})}{k}$$ (2.12)

In vector form, this scheme reads

$$\frac{U^j - U^{j-1}}{k} + A_h U^j = F^j \quad \text{for } j = 1, 2, 3, \dots, M + 1.$$ (2.13)

This scheme is called implicit, because the above formula is not a simple recurrence relation. Indeed, $U^{j+1}$ appears as the solution of an equation once $U^j$ is known. It is not a priori clear that this equation is solvable. In this particular case, we have

$$\begin{cases} U^j = (I + kA_h)^{-1}(U^{j-1} + kF^j) \quad \text{for } j = 1, 2, 3, \dots, M + 1 \\ U^0 = U_0 \end{cases}$$ (2.14)

16

since it is not hard to see that the matrix $I + kA_h$ is symmetric, positive definite, hence invertible.

In practical terms, the implementation of the backward Euler method entails the solution of a linear system at each time step, whereas the explicit method is simply a matrix-vector product and vector addition at each time step. The implicit method is thus more computationally intensive than the explicit method, but it has other benefits as we will see later.

The second example is the leapfrog or Richardson method, which is associated with the central differential quotient approximation of the time derivative

$$\frac{\partial u}{\partial t}(x_n, t_j) \approx \frac{u(x_n, t_{j+1}) - u(x_n, t_{j-1})}{2k} \tag{2.15}$$

Like the explicit Euler three point method we substitute this into the heat equation to get the equations for $U^{j+1}, U^j$ and $U^{j-1}$:

$$\frac{U^{j+1} - U^{j-1}}{2k} + A_h U^j = F^j \quad \text{for } j = 1, \dots, M. \tag{2.16}$$

Thus, we find that this method is an explicit two-step method since $U^{j+1}$ is explicitly given in terms of $U^j$ and $U^{j-1}$. Simplify the equation(2), we can get :

$$\begin{cases} U^{j+1} = U^{j-1} - 2kA_h U^j + 2kF^j \text{ for } j = 0, \cdots, M \\ U^0 = U_0 \ U^1 = U_1 \end{cases} \tag{2.17}$$

In particular, since this is a two-step method, we must somehow be ascribed to $U^1$ in order to initialize the recurrence, in addition to $U^1$.

## 2.2.3 Crank-Nicolson method

$$\begin{cases} \frac{u_i^{j+1} - u_i^j}{k} \quad + \frac{1}{2}\left(-\frac{u_{i+1}^{j+1} - 2u_i^{j+1} + u_{i-1}^{j+1}}{h^2} + \frac{u_{i+1}^{j+1} - u_{i-1}^{j+1}}{h}\right) \\ \quad + \frac{1}{2}\left(-\frac{u_{i+1}^j - 2u_i^j + u_{i-1}^j}{h^2} + \frac{u_{i+1}^j - u_{i-1}^j}{h}\right) = \frac{1}{2}f(x_i, t_{j+1}) + \frac{1}{2}f(x_i, t_j), \quad i \in \{1, \dots, N\}, j \in \{1, \dots, J\} \\ u_i^0 = 0, \qquad\qquad i \in \{1, \dots, N\} \\ u_0^j = u_{N+1}^j = 0, \qquad j \in \{0, \dots, J\} \end{cases} \tag{2.18}$$

From 2.18, we could establish directly the relationship

$$\frac{1}{k}(U^{j+1} - U^j) + \frac{1}{2}(\frac{1}{h^2}AU^{j+1} + \frac{1}{h}BU^{j+1}) + \frac{1}{2}(\frac{1}{h^2}AU^j + \frac{1}{h}BU^j) = \frac{1}{2}(F^{j+1} + F^j) \tag{2.19}$$

with

$$A = \begin{pmatrix} 2 & -1 & 0 & 0 & 0 & 0 & \dots & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 & \dots & 0 \\ 0 & 0 & -1 & 2 & -1 & 0 & \dots & 0 \\ 0 & 0 & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix} \quad B = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & -1 & 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & -1 & 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & -1 & 0 \end{pmatrix} \tag{2.20}$$

We reformulate the method 2.18 in a condensed way:

$$\begin{cases} \left(I + \dfrac{k}{2}C\right)U^{j+1} = \left(I - \dfrac{k}{2}C\right)U^j + \dfrac{k}{2}(F^{j+1} + F^j), \quad j \in \{1, \dots, J\} \\ U^0 = 0 \end{cases} \tag{2.21}$$

then $A, B, C$ satisfied $C = \dfrac{1}{h^2}A + \dfrac{1}{h}B$, we shall write $C$ explicitly:

$$C = \frac{1}{h^2}\begin{pmatrix} 2 & -1+h & 0 & 0 & 0 & \cdots & 0 \\ -1-h & 2 & -1+h & 0 & 0 & \cdots & 0 \\ 0 & -1-h & 2 & -1+h & 0 & \cdots & 0 \\ 0 & 0 & -1-h & 2 & -1+h & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 & 2 & -1+h \end{pmatrix} \tag{2.22}$$

This vector form 2.21 is finally realize in Python:

```python
def matrix_CF(N, func):
    # N as the number of knots (INCLUDING END POINTS)
    h = 1/(N-1)
    h_coeff = h**(-2)

    # Creation of diagonal part of C_h
    vect0 = h_coeff*2 * np.ones(N)
    C0 = np.diagflat(vect0, 0)

    # upper part of C_h
    vect1 = h_coeff*(-1+h) * np.ones(N-1)
    C1 = np.diagflat(vect1, 1)

    # lower part of C_h
    vect2 = h_coeff*(-1-h)*np.ones(N-1)
    C2 = np.diagflat(vect2, -1)

    # assembly C_h
    C_h = C0 + C1 + C2

    x = np.linspace(0, 1, N)
    F_h = func(x)

    return C_h, F_h
```

Listing 1: Martix representation of $C$ in python

Given $N, M \in \mathbb{N}^*$ the number of knots incluing end points, we take space step size $h = 1/(N-1)$ and set inner point $x_i = ih$, similally temporal step size is set as $k = 1/(M-1)$. The unknown values are $u_1^j, u_2^j, \dots, u_N^j$. For $i = 1, \cdots, N$ et $j = 1, \cdots, M+1$, the corresponding unknown vector is $U^j = (u_1^j, \dots, u_N^j)^T$. Because $u_0^j, u_1^j$ are in $\partial]0, 1[$, which is not defined here.

At each mesh point establishes the approximation $u_i^j \approx u(x_i, t_j)$ where in our case, the null frontier permets us to write $F^j$ as $(f_1^j, f_2^j, \cdots f_N^j)^T$.

It's recommended to use the sparse matrix to represent the matrix $C$, which could be implemented by `import scipy.sparse`. The core Python snippets to solve linear system 2.21 is

```python
def FDM_CN(N, M, f):
    h = 1/(N-1)
    k = 1/(M-1)
    t = np.linspace(0, 1, M)
    x = np.linspace(0, 1, N)
    U = np.zeros((N, M+1))
    # U[:, 0] = u(x,0) = 0 already satisfied

    C_h, F_h = matrix_CF(N, f)

```
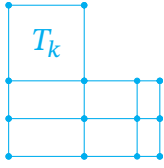
```
11      U[:, 0] = N*[0]
12      U[0, :] = U[-1, :] = (M+1)*[0]
13      for j in range(M):
14          U[:, j+1] = np.linalg.solve(np.eye(N)+k/2*C_h,\
15              np.matmul((np.eye(N) - k/2*C_h), U[:, j]) + k*F_h)
16
17      return t, U
```
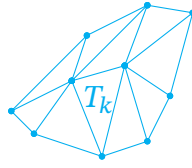
Listing 2: the Crank-Nicolson method for system linear system 2.21
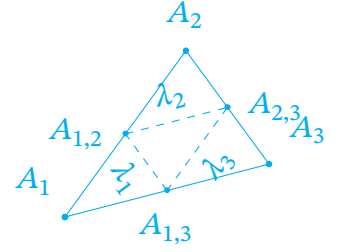
## 2.3 Finite Element Method (FEM) Approximation for 1D Heat Equation



(a) Rectangular $Q_1$ Finite Elements



(b) Triangular $P_1$ Lagrange Elements



(c) Triangular $P_2$ Lagrange Elements

Figure 5: The different FEM schemes

## 2.4 FDM for 2D Heat Equation

## 2.5 FEM for 2D Heat Equation

# 3 Analysis of Algorithms

## 3.1 Consistency

✠ **Definition 3.1** (truncation error)**:** *In a word, truncation error is an error caused by approximating a mathematical process.*

In Fourier analysis we have our classical truncation error (inducted from Bessel's inequality and Parseval's equality) defined by

$$\varepsilon_{\text{Fourier}}[M] = \left\| f - \widetilde{f}_M \right\|^2 = \sum_{k > M/2} \left| \left\langle f(\xi), e^{2i\pi k\xi} \right\rangle \right|^2 \tag{3.1}$$

In general FDM,

$$\varepsilon_{\text{FDM}} = \varepsilon_{h,k}(u)^j = \sum_{i=-m}^{l} B_i S_h(u_{t_{j+i}}) + \widetilde{F}_j, \quad m \le j \le \left( \frac{T}{k} - l \right) \tag{3.2}$$

✠ **Definition 3.2** (consistency)**:**

19

For different methods, we shall develop a general frame to

✠ **Definition 3.3:** *We say that the scheme (?) is consistent for the family of norms $\|\cdot\|_N$ if*

$$\max_{m \leq j \leq (\frac{T}{k}-l)} \left\| \varepsilon_{h,k}(u)^j \right\|_N \xrightarrow{(h,k) \to (0,0)} 0 \tag{3.3}$$

*And we say it's of order $p$ in space and $q$ in time for the family of norms $\|\cdot\|_N$ if*

$$\max_{m \leq j \leq (\frac{T}{k}-l)} \left\| \varepsilon_{h,k}(u)^j \right\|_N \leq C_u(h^p + k^q) \tag{3.4}$$

*Where $C_u$ is a constant that depends only on $u$.*

Consistency means that the scheme is trying its best to locally approximate the right numerical problem in the norm $\|\cdot\|_N$.

## 3.2  Stability

## 3.3  Order of Convergence

## 3.4  Possibility of Improvement

# 4  Application on Economics and Finance

## 4.1  The Black-Scholes PDE for Option Pricing

## 4.2  Optimal Portfolio for Consumption and Investment

# References

[Bre11]    Haim Brezis. *Functional Analysis, Sobolev Spaces and Partial Differential Equations*. Springer New York, 2011.

[Çin11]    Erhan Çinlar. *Probability and Stochastics*. Springer New York, Feb. 2011. 558 pp.

[Eva10]    Lawrence Evans. *Partial Differential Equations*. 2nd ed. American Mathematical Society, Mar. 2010.

[Fol99]    Gerald B. Folland. *Real Analysis*. 2nd ed. John Wiley & Sons, Mar. 1999. 406 pp.

[Gos20]    François Gosle. *Distributions, analyse de Fourier, équations aux dérivées partielles*. Ecole Polytechnique, Dec. 2020.

[Hab12]    Richard Haberman. *Applied partial differential equations with fourier series and boundary value problems*. 5th ed. Addison-Wesley, 2012.

[Hul17]    John C. Hull. *Options, Futures, and Other Derivatives*. 10th ed. PEARSON, Jan. 2017. 896 pp.

[Lax02]    Peter D. Lax. *Functional Analysis*. John Wiley & Sons, Mar. 2002. 604 pp.

[Lax07]    Peter D. Lax. *Linear Algebra and Its Applications*. 2nd ed. John Wiley & Sons, Aug. 2007. 394 pp.

[Luc16]    Hervé Le Dret; Brigitte Lucquin. *Partial Differential Equations: Modeling, Analysis and Numerical Approximation*. Springer International Publishing, 2016.

[Mal08]    Stephane Mallat. *A Wavelet Tour of Signal Processing*. 3rd ed. Elsevier Science Publishing Co Inc, Dec. 2008. 832 pp.

[Nef00]    Salih N. Neftci. *An Introduction to the Mathematics of Financial Derivatives*. 3rd ed. Elsevier Science & Techn., June 2000. 527 pp.

[Sch01]    Michelle Schatzman. *Analyse numérique : une approche mathématique*. 2nd ed. Paris: Dunod, 2001.

[Sha03]    Elias M. Stein; Rami Shakarchi. *Fourier Analysis: An Introduction*. PRINCETON UNIV PR, Apr. 2003. 328 pp.

[Shr10]    Steven Shreve. *Stochastic Calculus for Finance II*. Springer New York, Dec. 2010. 572 pp.

[Shr98]    Ioannis Karatzas; Steven E. Shreve. *Brownian Motion and Stochastic Calculus*. 2nd ed. Springer New York, 1998.

[Zil21]    Matthieu Bonnivard; Adina Ciomaga; Alessandro Zilio. "Méthodes numériques pour les EDO et les EDP". Notes de cours M1 Mathématiques. 2021.