Mathematical aspects of electronic structure theory

Eric CANCES

Ecole des Ponts and INRIA, Paris, France

Aussois, June 17-30, 2017

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x =unknown)

2

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x = unknown)

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$
 Solution: $\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x = unknown)

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$
 Solution: $\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

2

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x =unknown)

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$
Solution: $\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$
Solution: $\mathbf{x} = \begin{pmatrix} 1.082 \\ 0.862 \\ 1.035 \\ 0.979 \end{pmatrix}$

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x = unknown)

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1.082 \\ 0.862 \\ 1.035 \\ 0.979 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7.021 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x = unknown)

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1.082 \\ 0.862 \\ 1.035 \\ 0.979 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7.021 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} -2.77... \\ 7.19... \\ -0.51... \\ 1.90... \end{pmatrix}$$

Question 1: linear systems Ax = b ($A \in \mathbb{R}^{d \times d}$, $b \in \mathbb{R}^d$ given; x = unknown)

A computational code for solving Ax = b gives the following results.

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} 1.082 \\ 0.862 \\ 1.035 \\ 0.979 \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} 10 & 7.021 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} 32\\23\\33\\31 \end{pmatrix}$$

Solution:
$$\mathbf{x} = \begin{pmatrix} -2.77... \\ 7.19... \\ -0.51... \\ 1.90... \end{pmatrix}$$

Should you trust this code?

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Example: solve $\inf_{g(x)=0} E(x)$ where $E:\mathbb{R}^d \to \mathbb{R}$ and $g:\mathbb{R}^d \to \mathbb{R}^m$ are regular.

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Example: solve $\inf_{g(x)=0} E(x)$ where $E:\mathbb{R}^d \to \mathbb{R}$ and $g:\mathbb{R}^d \to \mathbb{R}^m$ are regular.

Introduce the Lagrangian $L:\mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ defined as

$$L(x,\lambda) = E(x) + \lambda^{T} g(x).$$

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Example: solve $\inf_{g(x)=0} E(x)$ where $E: \mathbb{R}^d \to \mathbb{R}$ and $g: \mathbb{R}^d \to \mathbb{R}^m$ are regular.

Introduce the Lagrangian $L:\mathbb{R}^d\times\mathbb{R}^m\to\mathbb{R}$ defined as

$$L(x,\lambda) = E(x) + \lambda^{T} g(x).$$

Then, the minimizers are obtained by solving the system of equations

$$\begin{cases} \nabla_x L(x,\lambda) = 0 \\ \nabla_\lambda L(x,\lambda) = 0, \end{cases}$$

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Example: solve $\inf_{g(x)=0} E(x)$ where $E:\mathbb{R}^d \to \mathbb{R}$ and $g:\mathbb{R}^d \to \mathbb{R}^m$ are regular.

Introduce the Lagrangian $L:\mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ defined as

$$L(x,\lambda) = E(x) + \lambda^{T} g(x).$$

Then, the minimizers are obtained by solving the system of equations

$$\begin{cases} \nabla_x L(x,\lambda) = 0 \\ \nabla_\lambda L(x,\lambda) = 0, \end{cases}$$

Application: d = 1, m = 1, E(x) = x, $g(x) = x^2$

Question 2: Lagrangian method for constrained optimization

Constrained optimization is ubiquitous in quantum physics and chemistry (e.g. Hartree-Fock, DFT, etc.). In Physics textbooks, such problems are solved using the Lagrangian method.

Example: solve $\inf_{g(x)=0} E(x)$ where $E: \mathbb{R}^d \to \mathbb{R}$ and $g: \mathbb{R}^d \to \mathbb{R}^m$ are regular.

Introduce the Lagrangian $L:\mathbb{R}^d\times\mathbb{R}^m\to\mathbb{R}$ defined as

$$L(x,\lambda) = E(x) + \lambda^{T} g(x).$$

Then, the minimizers are obtained by solving the system of equations

$$\begin{cases} \nabla_x L(x,\lambda) = 0 \\ \nabla_\lambda L(x,\lambda) = 0, \end{cases}$$

Application: d = 1, m = 1, E(x) = x, $g(x) = x^2$

$$\begin{cases} 1+2\lambda x=0 \\ x^2=0 \end{cases} \Rightarrow \begin{aligned} \text{No solution, though } x=0 \text{ is obviously a minimizer!} \\ \text{What's the catch?} \end{aligned}$$

Question 3: spectral theory of self-adjoint operators

Let $\mathcal H$ be a Hilbert space and H a self-adjoint operator on $\mathcal H$. In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of $\mathcal H$ such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Question 3: spectral theory of self-adjoint operators

Let \mathcal{H} be a Hilbert space and H a self-adjoint operator on \mathcal{H} . In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of \mathcal{H} such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Is this statement correct in general?

Question 3: spectral theory of self-adjoint operators

Let \mathcal{H} be a Hilbert space and H a self-adjoint operator on \mathcal{H} . In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of \mathcal{H} such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Is this statement correct in general?

Is is correct for a two-level quantum system?

Question 3: spectral theory of self-adjoint operators

Let $\mathcal H$ be a Hilbert space and H a self-adjoint operator on $\mathcal H$. In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of $\mathcal H$ such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Is this statement correct in general?

Is is correct for a two-level quantum system?

Is it correct for the harmonic oscillator?

Question 3: spectral theory of self-adjoint operators

Let \mathcal{H} be a Hilbert space and H a self-adjoint operator on \mathcal{H} . In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of \mathcal{H} such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Is this statement correct in general?

Is is correct for a two-level quantum system?

Is it correct for the harmonic oscillator?

Is it correct for the Schrödinger Hamiltonian of the hydrogen atom?

Question 3: spectral theory of self-adjoint operators

Let $\mathcal H$ be a Hilbert space and H a self-adjoint operator on $\mathcal H$. In Physics textbooks, it is claimed that there exists an orthonormal basis $(|\psi_n\rangle)$ of $\mathcal H$ such that

$$\langle \psi_m | \psi_n \rangle = \delta_{m,n}, \quad H | \psi_n \rangle = E_n | \psi_n \rangle, \quad E_n \in \mathbb{R}.$$

Is this statement correct in general?

Is is correct for a two-level quantum system?

Is it correct for the harmonic oscillator?

Is it correct for the Schrödinger Hamiltonian of the hydrogen atom?

Is it correct for the free-particle Hamiltonian?

- 1. A bit of numerical analysis
- 2. Constrained optimization and Lagrange multipliers
- 3. Spectral theory of self-adjoint operators

The deterministic models used in quantum physics and chemistry give rise to

- ullet linear eigenvalue problems (N-body Schrödinger eq., LR-TDDFT, BSE, ...)
- constrained optimization problems (HF, DFT, MCSCF, ...)
- algebraic equations (CC, ...)
- time-dependent linear or nonlinear Schrödinger equations (RT-TDDFT, ...)

The deterministic models used in quantum physics and chemistry give rise to

- ullet linear eigenvalue problems (N-body Schrödinger eq., LR-TDDFT, BSE, ...)
- constrained optimization problems (HF, DFT, MCSCF, ...)
- algebraic equations (CC, ...)
- time-dependent linear or nonlinear Schrödinger equations (RT-TDDFT, ...)

Solving numerically all these problems eventually boils down to (cleverly!) performing numerical quadratures and matrix-vector products.

The deterministic models used in quantum physics and chemistry give rise to

- ullet linear eigenvalue problems (N-body Schrödinger eq., LR-TDDFT, BSE, ...)
- constrained optimization problems (HF, DFT, MCSCF, ...)
- algebraic equations (CC, ...)
- time-dependent linear or nonlinear Schrödinger equations (RT-TDDFT, ...)

Solving numerically all these problems eventually boils down to (cleverly!) performing numerical quadratures and matrix-vector products.

Example: let $F: \mathbb{R}^d \to \mathbb{R}^d$. A standard iterative algorithm to solve the equation $F(\mathbf{x}) = 0$ is the Newton algorithm:

 \mathbf{x}_k begin given, solve the linear system $F'(\mathbf{x}_k)\mathbf{y}_k = -F(\mathbf{x}_k)$, then set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{y}_k$.

The deterministic models used in quantum physics and chemistry give rise to

- ullet linear eigenvalue problems (N-body Schrödinger eq., LR-TDDFT, BSE, ...)
- constrained optimization problems (HF, DFT, MCSCF, ...)
- algebraic equations (CC, ...)
- time-dependent linear or nonlinear Schrödinger equations (RT-TDDFT, ...)

Solving numerically all these problems eventually boils down to (cleverly!) performing numerical quadratures and matrix-vector products.

Example: let $F: \mathbb{R}^d \to \mathbb{R}^d$. A standard iterative algorithm to solve the equation $F(\mathbf{x}) = 0$ is the Newton algorithm:

 \mathbf{x}_k begin given, solve the linear system $F'(\mathbf{x}_k)\mathbf{y}_k = -F(\mathbf{x}_k)$, then set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{y}_k$.

Linear systems can themselves be solved by iterative algorithms based on matrix-vector products.

A key concept: conditioning

Consider a problem consisting of computing an output ${\bf s}$ from an input ${\bf y}$ (the data). The problem is called

- well-conditioned if a small variation of the input leads to a small variation of the output
- ill-conditioned otherwise.

A key concept: conditioning

Consider a problem consisting of computing an output ${\bf s}$ from an input ${\bf y}$ (the data). The problem is called

- well-conditioned if a small variation of the input leads to a small variation of the output
- ill-conditioned otherwise.

Toy example of a very ill-conditioned problem:

$$\mathbf{y} = \begin{pmatrix} 2 & 10^{17} \\ \mathbf{0} & 0.5 \end{pmatrix} \longrightarrow \mathbf{s} = \mathbf{eigenvalues of y} = (0.5; 2)$$

$$\mathbf{y} + \delta \mathbf{y} = \begin{pmatrix} 2 & 10^{17} \\ 10^{-17} & 0.5 \end{pmatrix} \longrightarrow \mathbf{s} + \delta \mathbf{s} = \mathbf{eigenvalues of y} + \delta \mathbf{y} = (0; 2.5).$$

An apparently nicer problem: solve the linear system Ax = b with

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \qquad \mathbf{and} \qquad \mathbf{b} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

The matrix A is symmetric, det(A) = 1, and

$$\mathbf{A}^{-1} = \begin{pmatrix} 25 & -41 & 10 & -6 \\ -41 & 68 & -17 & 10 \\ 10 & -17 & 5 & -3 \\ -6 & 10 & -3 & 2 \end{pmatrix}$$

Reference linear system

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} \\ \\ \\ \\ \\ \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

Solution =
$$\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

Slight perturbation of the right-hand side

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} \\ \\ \\ \\ \end{pmatrix} = \begin{pmatrix} 32.001 \\ 22.999 \\ 33.001 \\ 30.999 \end{pmatrix}$$

Solution =
$$\begin{pmatrix} 1.082 \\ 0.862 \\ 1.035 \\ 0.979 \end{pmatrix}$$

Slight modification of the matrix A

$$\begin{pmatrix} 10 & 7.021 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} \\ \\ \\ \\ \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}$$

Solution =
$$\begin{pmatrix} -2.77... \\ 7.19... \\ -0.51... \\ 1.90... \end{pmatrix}$$

This apparently nice problem is not so well-conditioned ...

$$\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$$
 for $1 \le p < +\infty$, $\|\mathbf{x}\|_{\infty} = \max_{1 \le i \le n} |x_i|$

$$\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$$
 for $1 \le p < +\infty$, $\|\mathbf{x}\|_{\infty} = \max_{1 \le i \le n} |x_i|$

 l^p -norm of a matrix $\mathbf{A} \in \mathbb{R}^{n imes m}$

$$\|\mathbf{A}\|_p := \sup_{\mathbf{x} \in \mathbb{R}^m \setminus \{0\}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}$$

$$\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$$
 for $1 \le p < +\infty$, $\|\mathbf{x}\|_{\infty} = \max_{1 \le i \le n} |x_i|$

 l^p -norm of a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$

$$\|\mathbf{A}\|_p := \sup_{\mathbf{x} \in \mathbb{R}^m \setminus \{0\}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}$$

Conditioning number: the conditioning number of the abstract problem $\mathbf{s} = f(\mathbf{y})$ at $\mathbf{y} = \mathbf{y}_0$ for the l^p -norm is $(\mathbf{s} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m)$ is

$$\kappa_p(\mathbf{y}_0) = \frac{\|f'(\mathbf{y}_0)\|_p \|\mathbf{y}_0\|_p}{\|f(\mathbf{y}_0)\|_p}.$$

$$\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$$
 for $1 \le p < +\infty$, $\|\mathbf{x}\|_{\infty} = \max_{1 \le i \le n} |x_i|$

 l^p -norm of a matrix $\mathbf{A} \in \mathbb{R}^{n imes m}$

$$\|\mathbf{A}\|_p := \sup_{\mathbf{x} \in \mathbb{R}^m \setminus \{0\}} \frac{\|\mathbf{A}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}$$

Conditioning number: the conditioning number of the abstract problem $\mathbf{s} = f(\mathbf{y})$ at $\mathbf{y} = \mathbf{y}_0$ for the l^p -norm is $(\mathbf{s} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m)$ is

$$\kappa_p(\mathbf{y}_0) = \frac{\|f'(\mathbf{y}_0)\|_p \|\mathbf{y}_0\|_p}{\|f(\mathbf{y}_0)\|_p}.$$

Rule of thumb: if the conditioning number is $\sim 10^p$ and if you compute in double precision ($\varepsilon_{\rm machine}=10^{-16}$), you can only trust the first 16-p digits of your result.

Conditioning number of an invertible square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ (for the l^p -norm)

$$\kappa_p(\mathbf{A}) := \|\mathbf{A}\|_p \|\mathbf{A}^{-1}\|_p$$

 $\kappa_p(\mathbf{A})$ is the max. w.r.t. x of the conditioning numbers of the problems:

- matrix-vector product: $y = (A, x) \mapsto s = Ax$
- linear system solver: $y = (A, x) \mapsto s = A^{-1}x$ (solve As = x)

Example:

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \longrightarrow \kappa_2(\mathbf{A}) = 2984 \quad \mathbf{and} \quad \kappa_{\infty}(\mathbf{A}) = 4488.$$

Theorem. Let $A \in \mathbb{R}^{n \times n}$ be an invertible matrix, and $b \in \mathbb{R}^n$, $b \neq 0$.

• Perturbation of the right-hand side

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \qquad \mathbf{A}(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b} \qquad \Rightarrow \qquad \frac{\|\delta \mathbf{x}\|_p}{\|\mathbf{x}\|_p} \le \kappa_p(\mathbf{A}) \frac{\|\delta \mathbf{b}\|_p}{\|\mathbf{b}\|_p}$$

and the inequality is optimal: A being given, there exists b and δb such that the inequality is an equality.

• Perturbation of the matrix

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \qquad (\mathbf{A} + \delta \mathbf{A}) (\mathbf{x} + \delta \mathbf{x}') = \mathbf{b} \qquad \Rightarrow \qquad \frac{\|\delta \mathbf{x}\|_p}{\|\mathbf{x} + \delta \mathbf{x}'\|_p} \le \kappa_p(\mathbf{A}) \frac{\|\delta \mathbf{A}\|_p}{\|\mathbf{A}\|_p}$$

and the inequality is optimal: A being given, there exists b and δA such that the inequality is an equality.

Properties of the conditioning number $\kappa_p(A)$

- $\kappa_p(\mathbf{A}) \geq 1$, $\forall \mathbf{A} \in \mathrm{GL}_n(\mathbb{R})$ (the set of invertible matrices)
- $\kappa_2(\mathbf{U}) = 1$ iff \mathbf{U} is orthogonal ($\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = I_n$)
- $1/\kappa_p(\mathbf{A})$ is a measure of the relative distance of the matrix \mathbf{A} to the set of singular matrices:

$$\frac{1}{\kappa_p(\mathbf{A})} = \min_{\mathbf{E} \mid (\mathbf{A} + \mathbf{E}) \notin GL_n(\mathbb{R})} \frac{\|\mathbf{E}\|_p}{\|\mathbf{A}\|_p}.$$

• If A is symmetric

$$\kappa_2(\mathbf{A}) = \frac{\max_i |\lambda_i(\mathbf{A})|}{\min_i |\lambda_i(\mathbf{A})|}$$

 $\lambda_1(A) \leq \lambda_2(A) \leq \cdots \leq \lambda_n(A)$ denoting the eigenvalues of A.

An iterative algorithm for solving a problem P is a method for constructing, from an initial guess x_0 , a sequence x_1 , x_2 , x_3 , ... such that (hopefully)

$$\mathbf{x}_k \xrightarrow[k \to +\infty]{} \mathbf{x},$$
 (1)

where x is a solution to the problem P (the solution if P is well-posed).

An iterative algorithm for solving a problem P is a method for constructing, from an initial guess x_0 , a sequence x_1 , x_2 , x_3 , ... such that (hopefully)

$$\mathbf{x}_k \xrightarrow[k \to +\infty]{} \mathbf{x},$$
 (1)

where x is a solution to the problem P (the solution if P is well-posed).

The algorithm is called convergent if (1) holds. In practice, the algorithm is stopped when some stopping criteria are met. The efficiency of the algorithm heavily relies on the choice of the stopping criteria.

An iterative algorithm for solving a problem P is a method for constructing, from an initial guess x_0 , a sequence x_1 , x_2 , x_3 , ... such that (hopefully)

$$\mathbf{x}_k \xrightarrow[k \to +\infty]{} \mathbf{x},$$
 (1)

where x is a solution to the problem P (the solution if P is well-posed).

The algorithm is called convergent if (1) holds. In practice, the algorithm is stopped when some stopping criteria are met. The efficiency of the algorithm heavily relies on the choice of the stopping criteria.

Examples of stopping test for linear systems Ax = b:

- a terrible one: maximum number of iterations $(k \ge k_{\text{max}}) \Rightarrow \text{STOP}$
- a good one: residual based error vector $(\|\mathbf{r}_k\|_2 \le \varepsilon_k) \Rightarrow$ STOP, where

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{A}(\mathbf{x} - \mathbf{x}_k), \quad \varepsilon_k = \varepsilon_{\text{tol}}(\|\mathbf{A}\|_1 \|\mathbf{x}_k\|_{\infty} + \|\mathbf{b}\|_2)$$
 (Oetli-Prager, 1963)

If A is symmetric, positive definite, then $\|\mathbf{r}_k\|_2 = \|\mathbf{x} - \mathbf{x}_k\|$ where $\|\cdot\|$ is the norm defined by $\|\mathbf{y}\| = \|\mathbf{A}\mathbf{y}\|_2$.

Reminder: gradient of a differentiable function $J: \mathbb{R}^d \to \mathbb{R}$

We have for all $\mathbf{x} \in \mathbb{R}^d$

$$\forall \mathbf{h} \in \mathbb{R}^d, \quad J(\mathbf{x} + \mathbf{h}) = J(\mathbf{x}) + \sum_{i=1}^d \frac{\partial J}{\partial x_i}(\mathbf{x}) \, h_i + o(\mathbf{h}) = J(\mathbf{x}) + \nabla J(\mathbf{x}) \cdot \mathbf{h} + o(\mathbf{h})$$

Euclidean scalar product

Euclidean gradient:
$$\nabla J(\mathbf{x}) = \begin{pmatrix} \frac{\partial J}{\partial x_1}(\mathbf{x}) \\ \cdot \\ \cdot \\ \frac{\partial J}{\partial x_d}(\mathbf{x}) \end{pmatrix}.$$

Reminder: gradient of a differentiable function $J: \mathbb{R}^d \to \mathbb{R}$

We have for all $\mathbf{x} \in \mathbb{R}^d$

$$\forall \mathbf{h} \in \mathbb{R}^d, \quad J(\mathbf{x} + \mathbf{h}) = J(\mathbf{x}) + \sum_{i=1}^d \frac{\partial J}{\partial x_i}(\mathbf{x}) \, h_i + o(\mathbf{h}) = J(\mathbf{x}) + \nabla J(\mathbf{x}) \cdot \mathbf{h} + o(\mathbf{h})$$

Euclidean scalar product

Euclidean gradient:
$$\nabla J(\mathbf{x}) = \begin{pmatrix} \frac{\partial J}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial J}{\partial x_d}(\mathbf{x}) \end{pmatrix}.$$

If \mathbb{R}^d is endowed with the scalar product $(\mathbf{x}, \mathbf{y})_S := \mathbf{x}^T S \mathbf{y}$, where $S \in \mathbb{R}^{d \times d}$ is a positive definite symmetric matrix, then the gradient of J, which we will denote by $\nabla_S J(\mathbf{x})$, is related to the Euclidean gradient $\nabla J(\mathbf{x})$ by

$$\nabla_S J(\mathbf{x}) = S^{-1} \nabla J(\mathbf{x}).$$

Geometrical interpretation of the gradient

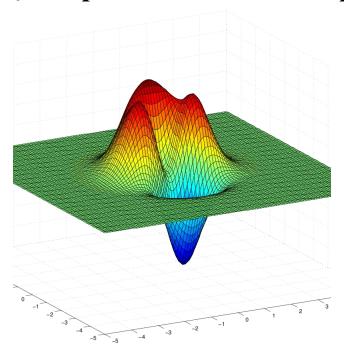
Let $J: \mathbb{R}^d \to \mathbb{R}$ of class C^1 , $\mathbf{x}_0 \in \mathbb{R}^d$ and $\alpha = J(\mathbf{x}_0)$. If $\nabla J(\mathbf{x}_0) \neq 0$, then

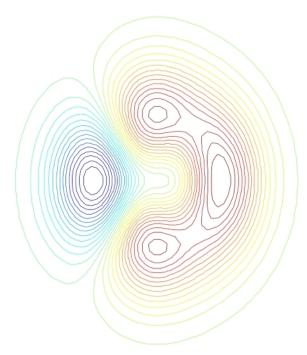
 \bullet in the vicinity of x_0 , the level set

$$\mathcal{C}_{\alpha} := \left\{ \mathbf{x} \in \mathbb{R}^d \mid J(\mathbf{x}) = \alpha \right\}$$

is a C^1 hypersurface (a codimension 1 C^1 manifold);

• the vector $\nabla J(\mathbf{x}_0)$ is orthogonal to the affine hyperplane tangent to \mathcal{C}_{α} at \mathbf{x}_0 and points toward the steepest ascent direction.





$$\mathbf{solve} \ \mathbf{A}\mathbf{x} = \mathbf{b} \qquad \Leftrightarrow \qquad \mathbf{solve} \ \min_{\mathbf{y} \in \mathbb{R}^d} J(\mathbf{y}) \quad \mathbf{where} \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}.$$

$$\mathbf{solve} \ \mathbf{A}\mathbf{x} = \mathbf{b} \qquad \Leftrightarrow \qquad \mathbf{solve} \ \min_{\mathbf{y} \in \mathbb{R}^d} J(\mathbf{y}) \quad \mathbf{where} \quad J(\mathbf{y}) := \frac{1}{2}\mathbf{y}^T \mathbf{A}\mathbf{y} - \mathbf{b}^T \mathbf{y}.$$

Gradient methods consist in choosing an initial guess $\mathbf{x}_0 \in \mathbb{R}^n$ and in building a sequence of iterates $(\mathbf{x}_k)_{k \in \mathbb{N}}$ of \mathbb{R}^n such that

$$J(\mathbf{x}_k) \underset{k \to +\infty}{\downarrow} \min_{\mathbb{R}^n} J$$
 Note that $\nabla J(\mathbf{y}) = \mathbf{A}\mathbf{y} - \mathbf{b}$

$$\mathbf{solve} \ \mathbf{A}\mathbf{x} = \mathbf{b} \qquad \Leftrightarrow \qquad \mathbf{solve} \ \min_{\mathbf{y} \in \mathbb{R}^d} J(\mathbf{y}) \quad \mathbf{where} \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}.$$

Gradient methods consist in choosing an initial guess $\mathbf{x}_0 \in \mathbb{R}^n$ and in building a sequence of iterates $(\mathbf{x}_k)_{k \in \mathbb{N}}$ of \mathbb{R}^n such that

$$J(\mathbf{x}_k) \underset{k \to +\infty}{\downarrow} \min_{\mathbb{R}^n} J$$
 Note that $\nabla J(\mathbf{y}) = \mathbf{A}\mathbf{y} - \mathbf{b}$

Gradient methods only involve matrix-vector and scalar products. There are particularly efficient when

- the matrix A cannot be stored (e.g. grid methods for Kohn-Sham)
- and/or matrix-vector products can be efficiently computed (sparse matrices, fast transforms such as FFT, ...)

$$\mathbf{solve} \ \mathbf{A}\mathbf{x} = \mathbf{b} \qquad \Leftrightarrow \qquad \mathbf{solve} \ \min_{\mathbf{y} \in \mathbb{R}^d} J(\mathbf{y}) \quad \mathbf{where} \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}.$$

Gradient methods consist in choosing an initial guess $\mathbf{x}_0 \in \mathbb{R}^n$ and in building a sequence of iterates $(\mathbf{x}_k)_{k \in \mathbb{N}}$ of \mathbb{R}^n such that

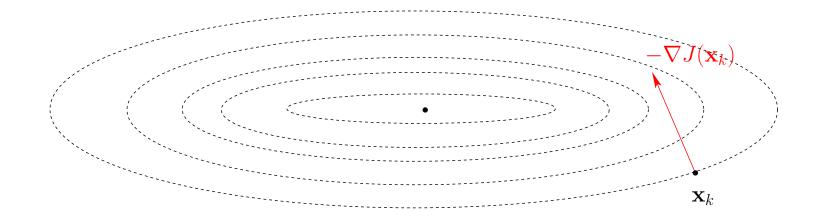
$$J(\mathbf{x}_k) \underset{k \to +\infty}{\downarrow} \min_{\mathbb{R}^n} J$$
 Note that $\nabla J(\mathbf{y}) = \mathbf{A}\mathbf{y} - \mathbf{b}$

Gradient methods only involve matrix-vector and scalar products. There are particularly efficient when

- the matrix A cannot be stored (e.g. grid methods for Kohn-Sham)
- and/or matrix-vector products can be efficiently computed (sparse matrices, fast transforms such as FFT, ...)

Remark: Extensions of gradient algorithms to general linear systems are available (MINRES - GMRES, 1986 - BiCGstab, 1992 - ...).

Fixed-step and optimal step gradient algorithms



The function J is decreasing in the direction

$$\mathbf{d}_k = -\nabla J(\mathbf{x}_k) = \mathbf{b} - \mathbf{A}\mathbf{x}_k$$
 (residual)

One then may choose

$$\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$$

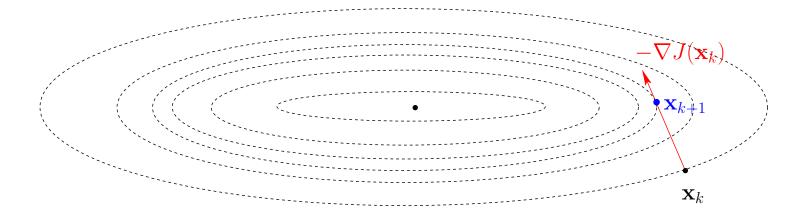
for some $t_k > 0$.

Fixed step: the step t is chosen once and for all

$$\left\{ egin{aligned} \mathbf{r}_k &= \mathbf{b} - \mathbf{A}\mathbf{x}_k \ \mathbf{x}_{k+1} &= \mathbf{x}_k + t\mathbf{r}_k \end{aligned}
ight.$$

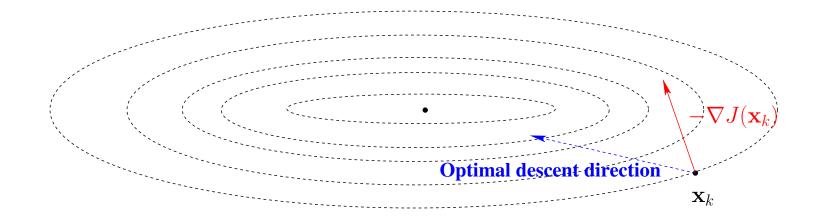
Optimal step: one chooses the "best" \mathbf{x}_{k+1} on the half-line $\mathbf{x}_k - t\nabla J(x_k)$

$$\begin{cases} \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k \\ t_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k} \\ \mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{r}_k \end{cases}$$



Conjugate gradient algorithm (1952)

The descent direction $\mathbf{d}_k = -\nabla J(\mathbf{x}_k)$ is optimal for infinitesimal steps, but not in general for finite step.



The conjugate gradient algorithm provides better descent directions d_k .

Conjugate gradient algorithm:

- Initialization. Choose $\mathbf{x}_0 \in \mathbb{R}^n$ and ε_{tol} , compute $\mathbf{r}_0 = \mathbf{b} \mathbf{A}\mathbf{x}_0$ and set $\mathbf{d}_0 = \mathbf{r}_0$. Set k = 0.
- Iterations.
 - 1. Stopping test: if $\|\mathbf{r}_k\|_2 \leq \varepsilon_{\text{tol}}(\|\mathbf{A}\|_1 \|\mathbf{x}_k\|_{\infty} + \|\mathbf{b}\|_2)$, stop.
 - **2.** Update x_k and the residual r_k :

$$\mathbf{z}_k = \mathbf{A}\mathbf{d}_k, \qquad \qquad t_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{d}_k^T \mathbf{z}_k}, \ \mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k, \qquad \qquad \mathbf{r}_{k+1} = \mathbf{r}_k - t_k \mathbf{z}_k,$$

3. Update the descent direction d_k :

$$eta_k = rac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k}, \qquad \mathbf{d}_{k+1} = \mathbf{r}_{k+1} + eta_k \mathbf{d}_k.$$

4. Set k = k + 1 and go to step **1.**

Krylov subspaces

The Krylov subspaces $(\mathcal{K}_k(\mathbf{y}))$ associated with a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and a vector \mathbf{y} are defined by

$$\mathcal{K}_k(\mathbf{y}) = \mathbf{Span}(\mathbf{y}, \mathbf{A}\mathbf{y}, \cdots, \mathbf{A}^k\mathbf{y})$$

Application to linear systems

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

$$= \mathbf{A}^{-1}(\mathbf{A}\mathbf{x}_0 + \mathbf{b} - \mathbf{A}\mathbf{x}_0)$$

$$= \mathbf{x}_0 + A^{-1}\mathbf{r}_0$$

$$= \mathbf{x}_0 + Q(\mathbf{A})\mathbf{r}_0 \quad \text{with } Q \text{ polynomial of degree } m \le n-1 \text{ (Hamilton-Cayley)}$$

$$\in \mathbf{x}_0 + \mathcal{K}_m(\mathbf{r}_0).$$

Theorem. Let (\mathbf{x}_k) the sequence generated by the conjugate gradient algorithm (with $\varepsilon_{\mathrm{tol}}=0$).

1. For all $k \ge 0$,

$$\mathbf{x}_k = \underset{\mathbf{y} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{r}_0)}{\operatorname{arginf}} J(\mathbf{y}), \qquad J(\mathbf{y}) = \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}$$

Theorem. Let (\mathbf{x}_k) the sequence generated by the conjugate gradient algorithm (with $\varepsilon_{\text{tol}} = 0$).

1. For all $k \ge 0$,

$$\mathbf{x}_k = \underset{\mathbf{y} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{r}_0)}{\operatorname{arginf}} J(\mathbf{y}), \qquad J(\mathbf{y}) = \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}$$

2. The sequence of Krylov subspace $\mathcal{K}_k(\mathbf{r}_0)$ is strictly increasing until the algorithm has converged: if $\mathbf{x}_k \neq \mathbf{x}$, dim $\mathcal{K}_k(\mathbf{r}_0) = k+1$. Consequently, the conjugate gradient algorithm converges in at most n iterations

Theorem. Let (\mathbf{x}_k) the sequence generated by the conjugate gradient algorithm (with $\varepsilon_{\text{tol}} = 0$).

1. For all $k \ge 0$,

$$\mathbf{x}_k = \underset{\mathbf{y} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{r}_0)}{\operatorname{arginf}} J(\mathbf{y}), \qquad J(\mathbf{y}) = \frac{1}{2} \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{b}^T \mathbf{y}$$

- 2. The sequence of Krylov subspace $\mathcal{K}_k(\mathbf{r}_0)$ is strictly increasing until the algorithm has converged: if $\mathbf{x}_k \neq \mathbf{x}$, dim $\mathcal{K}_k(\mathbf{r}_0) = k+1$. Consequently, the conjugate gradient algorithm converges in at most n iterations
- 3. If the conjugate gradient algorithm converges in m iterations, then $\forall 0 \le k \le m-1$,
 - $(\mathbf{r}_0, \mathbf{r}_1, \cdots, \mathbf{r}_k)$ is an orthogonal basis of $\mathcal{K}_k(\mathbf{r}_0)$: $\mathbf{r}_i^T \mathbf{r}_j = \delta_{ij}$
 - $(\mathbf{d}_0, \mathbf{d}_1, \cdots, \mathbf{d}_k)$ is an A-orthogonal basis of $\mathcal{K}_k(\mathbf{r}_0)$: $\mathbf{d}_i^T \mathbf{A} \mathbf{d}_j = \delta_{ij}$
 - \longrightarrow The descent directions d_k are A-conjugate

The conjugate gradient algorithm converges at least linearly

$$\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}} \le \rho^k \|\mathbf{x}_0 - \mathbf{x}\|_{\mathbf{A}}$$
 with $0 \le \rho = \left(\frac{\sqrt{\kappa_2(\mathbf{A})} - 1}{\sqrt{\kappa_2(\mathbf{A})} + 1}\right) < 1$,

where $\kappa_2(\mathbf{A}) = \frac{\lambda_n(\mathbf{A})}{\lambda_1(\mathbf{A})} \ge 1$ is the conditioning number of \mathbf{A} for the l^2 -norm, and where $\|\cdot\|_{\mathbf{A}}$ is the energy norm on \mathbb{R}^n defined by $\|\mathbf{y}\|_{\mathbf{A}} = (\mathbf{A}\mathbf{y},\mathbf{y})^{1/2}$.

The conjugate gradient algorithm converges at least linearly

$$\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}} \le \rho^k \|\mathbf{x}_0 - \mathbf{x}\|_{\mathbf{A}}$$
 with $0 \le \rho = \left(\frac{\sqrt{\kappa_2(\mathbf{A})} - 1}{\sqrt{\kappa_2(\mathbf{A})} + 1}\right) < 1$,

where $\kappa_2(\mathbf{A}) = \frac{\lambda_n(\mathbf{A})}{\lambda_1(\mathbf{A})} \ge 1$ is the conditioning number of \mathbf{A} for the l^2 -norm, and where $\|\cdot\|_{\mathbf{A}}$ is the energy norm on \mathbb{R}^n defined by $\|\mathbf{y}\|_{\mathbf{A}} = (\mathbf{A}\mathbf{y}, \mathbf{y})^{1/2}$.

Remarks

ullet This estimate is not optimal (convergence in at most n iterations)

The conjugate gradient algorithm converges at least linearly

$$\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}} \le \rho^k \|\mathbf{x}_0 - \mathbf{x}\|_{\mathbf{A}}$$
 with $0 \le \rho = \left(\frac{\sqrt{\kappa_2(\mathbf{A})} - 1}{\sqrt{\kappa_2(\mathbf{A})} + 1}\right) < 1$,

where $\kappa_2(\mathbf{A}) = \frac{\lambda_n(\mathbf{A})}{\lambda_1(\mathbf{A})} \ge 1$ is the conditioning number of \mathbf{A} for the l^2 -norm, and where $\|\cdot\|_{\mathbf{A}}$ is the energy norm on \mathbb{R}^n defined by $\|\mathbf{y}\|_{\mathbf{A}} = (\mathbf{A}\mathbf{y},\mathbf{y})^{1/2}$.

Remarks

- ullet This estimate is not optimal (convergence in at most n iterations)
- ullet The actual performance of the CG algorithm depends on the distribution of the eigenvalues of A

The conjugate gradient algorithm converges at least linearly

$$\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}} \le \rho^k \|\mathbf{x}_0 - \mathbf{x}\|_{\mathbf{A}}$$
 with $0 \le \rho = \left(\frac{\sqrt{\kappa_2(\mathbf{A})} - 1}{\sqrt{\kappa_2(\mathbf{A})} + 1}\right) < 1$,

where $\kappa_2(\mathbf{A}) = \frac{\lambda_n(\mathbf{A})}{\lambda_1(\mathbf{A})} \ge 1$ is the conditioning number of \mathbf{A} for the l^2 -norm, and where $\|\cdot\|_{\mathbf{A}}$ is the energy norm on \mathbb{R}^n defined by $\|\mathbf{y}\|_{\mathbf{A}} = (\mathbf{A}\mathbf{y},\mathbf{y})^{1/2}$.

Remarks

- ullet This estimate is not optimal (convergence in at most n iterations)
- ullet The actual performance of the CG algorithm depends on the distribution of the eigenvalues of A
- The smaller the conditioning number, the faster the algorithm

The conjugate gradient algorithm converges at least linearly

$$\|\mathbf{x}_k - \mathbf{x}\|_{\mathbf{A}} \le \rho^k \|\mathbf{x}_0 - \mathbf{x}\|_{\mathbf{A}}$$
 with $0 \le \rho = \left(\frac{\sqrt{\kappa_2(\mathbf{A})} - 1}{\sqrt{\kappa_2(\mathbf{A})} + 1}\right) < 1$,

where $\kappa_2(\mathbf{A}) = \frac{\lambda_n(\mathbf{A})}{\lambda_1(\mathbf{A})} \ge 1$ is the conditioning number of \mathbf{A} for the l^2 -norm, and where $\|\cdot\|_{\mathbf{A}}$ is the energy norm on \mathbb{R}^n defined by $\|\mathbf{y}\|_{\mathbf{A}} = (\mathbf{A}\mathbf{y},\mathbf{y})^{1/2}$.

Remarks

- ullet This estimate is not optimal (convergence in at most n iterations)
- ullet The actual performance of the CG algorithm depends on the distribution of the eigenvalues of A
- The smaller the conditioning number, the faster the algorithm
- --> Preconditioning can (often must) be used to reduced the cond. numb.

Iterative algorithms are usually totally inefficient without preconditioning.

Preconditioning of linear systems:

Basic idea: instead of solving

$$Ax = b$$

solve

$$\left\{ \begin{array}{l} \mathbf{P}^{-1/2}\mathbf{A}\mathbf{P}^{-1/2}\mathbf{z} = \mathbf{P}^{-1/2}\mathbf{b}, \\ \mathbf{P}^{1/2}\mathbf{x} = \mathbf{z}. \end{array} \right.$$

for some symmetric matrix P such that

$$\kappa_2(\mathbf{P}^{-1/2}\mathbf{A}\mathbf{P}^{-1/2}) \ll \kappa_2(\mathbf{A})$$

This replacement can be done implicitely: no need to compute $P^{-1/2}$.

Preconditioned conjugate gradient algorithm

- Initialisation. Choose $\mathbf{x}_0 \in \mathbb{R}^n$ and a threshold $\varepsilon_{\mathrm{tol}}$, compute $\mathbf{r}_0 = \mathbf{b} \mathbf{A}\mathbf{x}_0$, and the solution \mathbf{y}_0 to $\mathbf{P}\mathbf{y}_0 = \mathbf{r}_0$. Set $\mathbf{d}_0 = \mathbf{y}_0$ and k = 0.;
- Iterations.
 - **1. Stopping test: if** $\|\mathbf{r}_{k}\|_{2} \leq \varepsilon_{\text{tol}}(\|A\|_{1}\|x_{k}\|_{\infty} + \|b\|_{2})$, stop.
 - **2.** Update x_k and r_k

$$\mathbf{z}_k = \mathbf{A}\mathbf{d}_k, \qquad \qquad t_k = rac{\mathbf{y}_k^T \mathbf{r}_k}{\mathbf{d}_k^T \mathbf{z}_k}, \ \mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k, \qquad \qquad \mathbf{r}_{k+1} = \mathbf{r}_k - t_k \mathbf{z}_k, \ \mathbf{Solve} \quad \mathbf{P}\mathbf{y}_{k+1} = \mathbf{r}_{k+1}$$

3. Updated the descent direction d_k

$$eta_k = rac{\mathbf{y}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{y}_k^T \mathbf{r}_k}, \qquad \mathbf{d}_{k+1} = \mathbf{y}_{k+1} + eta_k \mathbf{d}_k.$$

4. Set k = k + 1 and go to step 1.

For the preconditioning technique to be efficient, the preconditioner ${\cal P}$ must fulfill two conditions

1.
$$\kappa_2(\mathbf{P}^{-1/2}\mathbf{A}\mathbf{P}^{-1/2}) \ll \kappa_2(\mathbf{A})$$

2. linear systems of the form Py = r are easy to solve.

 \longrightarrow A trade-off has to be made.

- "Algebraic preconditioners"
 - diagonal preconditioner
 - SSOR preconditioner
 - incomplete LU or Cholesky decomposition
- "Physical preconditioners"
 - Multigrid methods
 - Simplified model

Example: planewave discretization of periodic Schrödinger operators

$$H = -\frac{1}{2}\frac{d^2}{dx^2} + V$$
, $V(x) = |\cos(\pi x)|$, $e_k(x) = e^{2i\pi kx}$, $X_N = \mathbf{Span}(e_k, |k| \le N)$

$$H_{kl} = \langle e_k | H | e_l \rangle = 2\pi^2 |k|^2 \delta_{kl} + \hat{V}_{kl}, \quad \hat{V}_{kl} = \int_0^1 V(x) \, e^{2i\pi(l-k)x} \, dx, \quad -N \le k, l \le N$$

Solve
$$\mathbf{H}\mathbf{x} = \mathbf{b}$$
, with $\mathbf{b} = (1, \dots, 1)^T$

 \longrightarrow Possible preconditioner: P s.t. $P_{kl} = (1 + 2\pi^2 |k|^2) \delta_{kl}$

Stopping criterion: $\|\mathbf{r}_k\|_2 \leq 10^{-10}$ where $\mathbf{r}_k = \mathbf{b} - \mathbf{H}\mathbf{x}_k$

N	Size of the matrix H	# CG iter.	# PCG iter.
50	101	71	5
100	201	98	5
200	401	304	5
400	801	613	5

2 - Constrained optimization and Lagrange multipliers

Let $E:\mathbb{R}^d\to\mathbb{R}$ and $g:\mathbb{R}^d\to\mathbb{R}^m$ be two differentiable functions and consider the optimization problem

$$\inf_{\mathbf{x} \in K} E(\mathbf{x}) \quad \text{where} \quad K = \left\{ \mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) = 0 \right\}.$$

Let $E:\mathbb{R}^d\to\mathbb{R}$ and $g:\mathbb{R}^d\to\mathbb{R}^m$ be two differentiable functions and consider the optimization problem

$$\inf_{\mathbf{x} \in K} E(\mathbf{x}) \quad \text{where} \quad K = \left\{ \mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) = 0 \right\}.$$

Definition (qualification of the constraints). The equality constraints g=0 are called qualified at $\mathbf{x}_0 \in K$ if $g'(\mathbf{x}_0) \in \mathbb{R}^{m \times d}$ is surjective (i.e. $\operatorname{Ran}(g'(\mathbf{x}_0)) = \mathbb{R}^m$).

Let $E:\mathbb{R}^d\to\mathbb{R}$ and $g:\mathbb{R}^d\to\mathbb{R}^m$ be two differentiable functions and consider the optimization problem

$$\inf_{\mathbf{x} \in K} E(\mathbf{x}) \quad \text{where} \quad K = \left\{ \mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) = 0 \right\}.$$

Definition (qualification of the constraints). The equality constraints g=0 are called qualified at $\mathbf{x}_0 \in K$ if $g'(\mathbf{x}_0) \in \mathbb{R}^{m \times d}$ is surjective (i.e. $\operatorname{Ran}(g'(\mathbf{x}_0)) = \mathbb{R}^m$).

Theorem (Euler-Lagrange theorem). Let $x_0 \in K$ be a local minimum of E on K. Assume that

- 1. $\mathbf{x} \mapsto g'(\mathbf{x})$ is continuous in the vicinity of \mathbf{x}_0 ;
- **2.** the equality constraints g = 0 is qualified at \mathbf{x}_0 .

Then, there exists a unique $\lambda \in \mathbb{R}^m$ such that

$$\nabla E(\mathbf{x}_0) + g'(\mathbf{x}_0)^T \lambda = 0,$$

where $g'(\mathbf{x}_0)^T$ is the transpose of $g'(\mathbf{x}_0)$. The vector λ is called the Lagrange multiplier of the constraint g=0.

Euler-Lagrange equations

Assume that the constraints are qualified at any point of K. Then solving

$$\begin{cases} \mathbf{seek} \ (\mathbf{x}, \lambda) \in \mathbb{R}^d \times \mathbb{R}^m \ \mathbf{such that} \\ \nabla E(\mathbf{x}) + g'(\mathbf{x})^T \lambda = 0 \\ g(\mathbf{x}) = 0 \end{cases} \tag{2}$$

allows one to find all the critical points (among which the local minimizers and the local maximizers) of E on K.

Remark : the above problem consists of (d+m) scalar equations with (d+m) scalar unknowns.

Euler-Lagrange equations

Assume that the constraints are qualified at any point of K. Then solving

$$\begin{cases} \mathbf{seek} \ (\mathbf{x}, \lambda) \in \mathbb{R}^d \times \mathbb{R}^m \ \mathbf{such that} \\ \nabla E(\mathbf{x}) + g'(\mathbf{x})^T \lambda = 0 \\ g(\mathbf{x}) = 0 \end{cases} \tag{3}$$

allows one to find all the critical points (among which the local minimizers and the local maximizers) of E on K.

Remark : the above problem consists of (d+m) scalar equations with (d+m) scalar unknowns.

The solutions of the Euler-Lagrange equations (4) are called the critical points of E on K.

Euler-Lagrange equations

Assume that the constraints are qualified at any point of K. Then solving

$$\begin{cases} \mathbf{seek} \ (\mathbf{x}, \lambda) \in \mathbb{R}^d \times \mathbb{R}^m \ \mathbf{such that} \\ \nabla E(\mathbf{x}) + g'(\mathbf{x})^T \lambda = 0 \\ g(\mathbf{x}) = 0 \end{cases} \tag{4}$$

allows one to find all the critical points (among which the local minimizers and the local maximizers) of E on K.

Remark: the above problem consists of (d+m) scalar equations with (d+m) scalar unknowns.

The solutions of the Euler-Lagrange equations (4) are called the critical points of E on K.

Remark. Equations (4) are equivalent to seeking $(\mathbf{x}, \lambda) \in \mathbb{R}^d \times \mathbb{R}^m$ s.t.

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = 0, \ \, \nabla_{\lambda} L(\mathbf{x}, \lambda) = 0, \quad \text{where} \quad L(\mathbf{x}, \lambda) := E(\mathbf{x}) + \lambda \cdot g(\mathbf{x}) \quad \text{(Lagrangian)}.$$

Very important take-home messages

A mathematical theorem consists of

- a list of assumptions;
- one of more results following from these assumptions.

Do not forget to check the assumptions before using the results!

Very important take-home messages

A mathematical theorem consists of

- a list of assumptions;
- one of more results following from these assumptions.

Do not forget to check the assumptions before using the results!

Back to the example d=1, m=1, E(x)=x, $g(x)=x^2$. Then

$$K = \{x \in \mathbb{R} \mid g(x) = 0\} = \{0\}$$
 and $g'(0) = 0$.

The constraint g=0 is therefore not qualified, and this is the reason why the Lagragian method fails!

Very important take-home messages

A mathematical theorem consists of

- a list of assumptions;
- one of more results following from these assumptions.

Do not forget to check the assumptions before using the results!

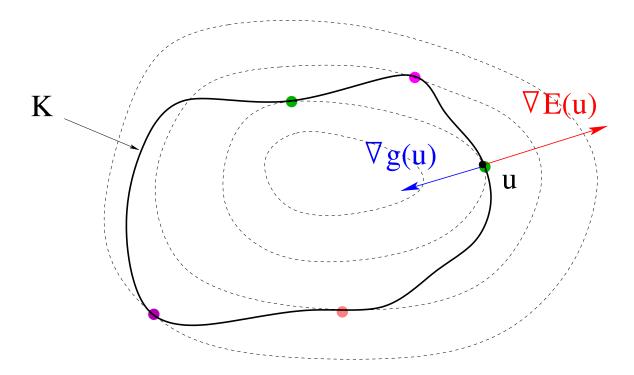
Back to the example
$$d = 1$$
, $m = 1$, $E(x) = x$, $g(x) = x^2$. Then $K = \{x \in \mathbb{R} \mid g(x) = 0\} = \{0\}$ and $g'(0) = 0$.

The constraint g=0 is therefore not qualified, and this is the reason why the Lagragian method fails!

Be all the more careful that not every "reasonable" mathematical statement is true!

Example: let \mathcal{H} be a Hilbert space. A continuous function $E:\mathcal{H}\to\mathbb{R}$ going to $+\infty$ at infinity does not necessarily have a minimizer.

A simple 2D example (d = 2, m = 1)



On
$$K = g^{-1}(0) = \{ \mathbf{x} \in \mathbb{R}^2 \mid g(\mathbf{x}) = 0 \}$$
, the function E possesses

- two local minimizers, all global
- two local maximizers, among which the global maximizer
- one critical point which is neither a local minimizer not a local maximizer.

Sketch of the proof

- Let \mathbf{x}_0 be a local minimizer of E on $K = g^{-1}(0) = \{\mathbf{x} \in \mathbb{R}^d \mid g(\mathbf{x}) = 0\}$ and $\alpha = E(\mathbf{x}_0)$.
- If the constraint g=0 is qualified at \mathbf{x}_0 (i.e. if $g'(\mathbf{x}_0):\mathcal{H}\to\mathcal{K}$ is surjective), then, in the vicinity of \mathbf{x}_0 , K is a C^1 manifold with tangent space

$$T_{\mathbf{x}_0}K = \left\{ \mathbf{h} \in \mathbb{R}^d \mid g'(\mathbf{x}_0)\mathbf{h} = 0 \right\} = \mathbf{Ker}(g'(\mathbf{x}_0)).$$

• Since \mathbf{x}_0 is a minimizer of E on K, the vector $\nabla E(\mathbf{x}_0)$ must be orthogonal to $T_{\mathbf{x}_0}K$. Indeed, for any $\mathbf{h} \in T_{\mathbf{x}_0}K$, there exists a C^1 curve $\phi: [-1,1] \to \mathbb{R}^d$ drawn on K such that $\phi(0) = \mathbf{x}_0$ et $\phi'(0) = \mathbf{h}$, and we have

$$0 \le E(\phi(t)) - E(\mathbf{x}_0) = E(\mathbf{x}_0 + t\mathbf{h} + o(t)) - E(\mathbf{x}_0) = t\nabla E(\mathbf{x}_0) \cdot \mathbf{h} + o(t).$$

• We have

$$\nabla E(\mathbf{x}_0) \in (T_{\mathbf{x}_0}K)^{\perp} = (\mathbf{Ker}(g'(\mathbf{x}_0)))^{\perp} = \mathbf{Ran}(g'(\mathbf{x}_0)^T).$$

ullet Therefore, there exists $\lambda \in \mathbb{R}^m$ such that $\nabla E(\mathbf{x}_0) + g'(\mathbf{x}_0)^T \lambda = 0$.

Remarks

- The above results can be extended to the case when $E:\mathcal{H}\to\mathbb{R}$ and $g:\mathcal{H}\to\mathcal{K}$ where \mathcal{H} and \mathcal{K} are Hilbert spaces.
- Most often, Lagrange multipliers have a "physical" interpretation
 - statistical mechanics, the equilibrium state of a chemical system interacting with its environment is obtained by maximizing the entropy under the constraints that the energy, the volume and the concentration of chemical species are given on average:
 - \rightarrow the Lagrange multipliers are respectively 1/T, P/T and μ_i/T (T: temperature, P: pressure, μ_i chemical potential of species i)
 - fluid mechanics, the admissible dynamics of an incompressible fluid are the critical points of the action under the constraint that the density of the fluid remains constant ($\operatorname{div}(u) = 0$)
 - \rightarrow the Lagrange multiplier of the incompressibility constraint is the pressure field.

Analytical derivatives

$$\forall \mathbf{R} \in \mathbb{R}^k, \quad W(\mathbf{R}) = \inf \left\{ E(\mathbf{R}, \mathbf{x}), \ \mathbf{x} \in \mathbb{R}^d, \ g(\mathbf{R}, \mathbf{x}) = 0 \right\}$$
 (5)

with $E: \mathbb{R}^k \times \mathbb{R}^d \to \mathbb{R}$, $g: \mathbb{R}^k \times \mathbb{R}^d \to \mathbb{R}^m$.

Assume (5) has a unique minimizer x(R) and $R \mapsto x(R)$ is regular. Then,

$$W(\mathbf{R}) = E(\mathbf{R}, \mathbf{x}(\mathbf{R})) \quad \Rightarrow \quad \frac{\partial W}{\partial R_k}(\mathbf{R}) = \frac{\partial E}{\partial R_k}(\mathbf{R}, \mathbf{x}(\mathbf{R})) + \nabla_{\mathbf{x}} E(\mathbf{R}, \mathbf{x}(\mathbf{R})) \cdot \frac{\partial \mathbf{x}}{\partial R_k}(\mathbf{R}),$$

$$g(\mathbf{R}, \mathbf{x}(\mathbf{R})) = 0 \implies \frac{\partial g}{\partial R_k}(\mathbf{R}, \mathbf{x}(\mathbf{R})) + g'_{\mathbf{x}}(\mathbf{R}, \mathbf{x}(\mathbf{R})) \frac{\partial \mathbf{x}}{\partial R_k}(\mathbf{R}) = 0.$$

Euler-Lagrange equation: $\nabla_{\mathbf{x}} E(\mathbf{R}, \mathbf{x}(\mathbf{R})) + g'_{\mathbf{x}}(\mathbf{R}, \mathbf{x}(\mathbf{R}))^T \lambda(\mathbf{R}) = 0.$

Therefore
$$\frac{\partial W}{\partial R_k}(\mathbf{R}) = \frac{\partial E}{\partial R_k}(\mathbf{R}, \mathbf{x}(\mathbf{R})) + \left(\frac{\partial g}{\partial R_k}(\mathbf{R}, \mathbf{x}(\mathbf{R})), \lambda(\mathbf{R})\right).$$

3 - Spectral theory of self-adjoint operators

References:

- E.B. Davies, *Linear operators and their spectra*, Cambridge University Press 2007.
- B. Helffer, Spectral theory and its applications, Cambridge University Press 2013.
- M. Reed and B. Simon, *Modern methods in mathematical physics*, in 4 volumes, 2nd edition, Academic Press 1972-1980.

Notation: in this section, \mathcal{H} denotes a separable complex Hilbert space, $\langle\cdot|\cdot\rangle$ its scalar product, and $\|\cdot\|$ the associated norm.

1. Each closed quantum system is associated with a separable complex Hilbert space $\mathcal{H}.$

- 1. Each closed quantum system is associated with a separable complex Hilbert space \mathcal{H} .
- 2. If the state of the system at time t is completely known, it can be described by a normalized vector $\psi(t)$ of \mathcal{H} . It is said to be in a pure state. The set of pure states is diffeomorphic to the projective space $P(\mathcal{H})$.

- 1. Each closed quantum system is associated with a separable complex Hilbert space \mathcal{H} .
- 2. If the state of the system at time t is completely known, it can be described by a normalized vector $\psi(t)$ of \mathcal{H} . It is said to be in a pure state. The set of pure states is diffeomorphic to the projective space $P(\mathcal{H})$.
- 3. Physical observables are represented by self-adjoint operators on \mathcal{H} .

- 1. Each closed quantum system is associated with a separable complex Hilbert space \mathcal{H} .
- 2. If the state of the system at time t is completely known, it can be described by a normalized vector $\psi(t)$ of \mathcal{H} . It is said to be in a pure state. The set of pure states is diffeomorphic to the projective space $P(\mathcal{H})$.
- 3. Physical observables are represented by self-adjoint operators on \mathcal{H} .
- 4. Let a be a physical observable represented by the self-adjoint operator A. The outcome of a measurement of a is always in $\sigma(A)$, the spectrum of A.

- 1. Each closed quantum system is associated with a separable complex Hilbert space \mathcal{H} .
- 2. If the state of the system at time t is completely known, it can be described by a normalized vector $\psi(t)$ of \mathcal{H} . It is said to be in a pure state. The set of pure states is diffeomorphic to the projective space $P(\mathcal{H})$.
- 3. Physical observables are represented by self-adjoint operators on \mathcal{H} .
- 4. Let a be a physical observable represented by the self-adjoint operator A. The outcome of a measurement of a is always in $\sigma(A)$, the spectrum of A.
- 5. If, just before the measurement, the system is in the pure state $\psi(t_0)$, then the probability that the outcome lays in the interval $B \subset \mathbb{R}$ is $\|\mathbb{1}_B(A)\psi(t_0)\|^2$, where $\mathbb{1}_B$ is the characteristic function of B and $\mathbb{1}_B(A)$ is defined by functional calculus.

- 1. Each closed quantum system is associated with a separable complex Hilbert space \mathcal{H} .
- 2. If the state of the system at time t is completely known, it can be described by a normalized vector $\psi(t)$ of \mathcal{H} . It is said to be in a pure state. The set of pure states is diffeomorphic to the projective space $P(\mathcal{H})$.
- 3. Physical observables are represented by self-adjoint operators on \mathcal{H} .
- 4. Let a be a physical observable represented by the self-adjoint operator A. The outcome of a measurement of a is always in $\sigma(A)$, the spectrum of A.
- 5. If, just before the measurement, the system is in the pure state $\psi(t_0)$, then the probability that the outcome lays in the interval $B \subset \mathbb{R}$ is $\|\mathbb{1}_B(A)\psi(t_0)\|^2$, where $\mathbb{1}_B$ is the characteristic function of B and $\mathbb{1}_B(A)$ is defined by functional calculus.
- 6. If the system is isolated, its dynamics between two successive measures is given by $\psi(t) = U(t-t_0)\psi(t_0)$ where $U(\tau) = e^{-i\tau H/\hbar}$, H being the Hamiltonian, i.e. the self-adjoint operator associated with the energy.

Definition (Hilbert space). A Hilbert space is a real or complex vector space \mathcal{H} endowed with a scalar product $\langle \cdot | \cdot \rangle$ and complete for the associated norm $\| \cdot \|$.

Definition (completeness). A sequence $(\psi_n)_{n\in\mathbb{N}}$ of elements of a normed vector space $(\mathcal{H},\|\cdot\|)$ is Cauchy if $\forall \varepsilon>0, \quad \exists N\in\mathbb{N} \quad \text{s.t.} \quad \forall q\geq p\geq N, \quad \|\psi_p-\psi_q\|\leq \varepsilon.$

The normed vector space $(\mathcal{H}, \|\cdot\|)$ is complete if any Cauchy sequence of elements of \mathcal{H} converges in \mathcal{H} .

Example: all finite dimensional normed \mathbb{R} - or \mathbb{C} -vector spaces are complete.

• Endowed with the Euclidean scalar product, \mathbb{C}^d is a Hilbert space:

$$(\mathbf{x}, \mathbf{y})_2 = \sum_{1 \le i \le d} \overline{x_i} y_i, \qquad \|\mathbf{x}\|_2 = (\mathbf{x}, \mathbf{x})_2^{1/2} = \left(\sum_{1 \le i \le d} |x_i|^2\right)^{1/2}.$$

ullet Let $S \in \mathbb{C}^{d \times d}$ be a positive definite hermitian matrix

$$(S_{ji} = \overline{S_{ij}} \text{ for all } 1 \leq i, j \leq d \text{ and } \mathbf{x}^* S \mathbf{x} > 0 \text{ for all } \mathbf{x} \in \mathbb{C}^d \setminus \{0\}).$$

Then $(\mathbf{x}, \mathbf{y})_S = \mathbf{x}^* S \mathbf{y}$ defines a scalar product on \mathbb{C}^d and

$$\forall \mathbf{x} \in \mathbb{C}^d$$
, $\lambda_1(S) \|\mathbf{x}\|_2 \le \|\mathbf{x}\|_S \le \lambda_d(S) \|\mathbf{x}\|_2$,

where $\lambda_1(S) \leq \lambda_2(S) \leq \cdots \leq \lambda_d(S)$ are the eigenvalues of S.

Fundamental examples: the Hilbert space $L^2(\mathbb{R}^d, \mathbb{C})$.

• The sequilinear form

$$(u,v) \mapsto (u,v)_{L^2} := \int_{\mathbb{R}^d} \overline{u}v := \int_{\mathbb{R}^d} \overline{u}(\mathbf{r}) v(\mathbf{r}) d\mathbf{r}$$

defines a scalar product on

$$C_{\mathrm{c}}^{\infty}(\mathbb{R}^d,\mathbb{C}):=\left\{v\in C^{\infty}(\mathbb{R}^d,\mathbb{C})\mid v=0 \text{ outside some bounded set}\right\},$$

but $C^\infty_{\rm c}(\mathbb{R}^d,\mathbb{C})$, endowed with the scalar product $(\cdot,\cdot)_{L^2}$, is not a Hilbert space.

• To obtain a Hilbert space, we have to "complete" it with "all the limits of the Cauchy sequences of elements of $C_c^{\infty}(\mathbb{R}^d)$ ". We thus obtain the set

$$L^{2}(\mathbb{R}^{d}, \mathbb{C}) := \left\{ u : \mathbb{R}^{d} \to \mathbb{C} \mid \int_{\mathbb{R}^{d}} |u|^{2} < \infty \right\},$$

which, endowed with the scalar product $(u, v)_{L^2}$, is a Hilbert space.

- Technical details:
 - one must use the Lebesgue integral (doesn't work with Riemann integral);
 - the elements of $L^2(\mathbb{R}^d,\mathbb{C})$ are in fact equivalence classes of measurable functions (for the Lebesgue measure) for the equivalence relation $u \sim v$ iff u = v everywhere except possibly on a set of Lebesgue measure equal to zero.

Fundamental examples: the Sobolev spaces $H^1(\mathbb{R}^d,\mathbb{C})$ and $H^2(\mathbb{R}^d,\mathbb{C})$.

• The sets

$$H^{1}(\mathbb{R}^{d},\mathbb{C}) := \left\{ u \in L^{2}(\mathbb{R}^{d},\mathbb{C}) \mid \nabla u \in (L^{2}(\mathbb{R}^{d},\mathbb{C}))^{d} \right\},$$

$$H^{2}(\mathbb{R}^{d},\mathbb{C}) := \left\{ u \in L^{2}(\mathbb{R}^{d},\mathbb{C}) \mid \nabla u \in (L^{2}(\mathbb{R}^{d},\mathbb{C}))^{d} \text{ and } D^{2}u \in (L^{2}(\mathbb{R}^{d},\mathbb{C}))^{d \times d} \right\}$$

are vector spaces. Respectively endowed with the scalar products

$$(u,v)_{H^1} := \int_{\mathbb{R}^d} \overline{u}v + \int_{\mathbb{R}^d} \overline{\nabla u} \cdot \nabla v,$$

$$(u,v)_{H^2} := \int_{\mathbb{R}^d} \overline{u}v + \int_{\mathbb{R}^d} \overline{\nabla u} \cdot \nabla v + \int_{\mathbb{R}^d} \overline{D^2u} : D^2v,$$

they are Hilbert spaces.

• Technical detail: the gradient and the second derivatives are defined by means of distribution theory.

Remark. Let $u \in H^1(\mathbb{R}^d)$. A function $\widetilde{u} \in H^1(\mathbb{R}^d)$ can be a very accurate approximation of u in $L^2(\mathbb{R}^d)$ and a terrible approximation of u in $H^1(\mathbb{R}^d)$.

For instance, let
$$u(x)=\frac{1}{1+x^2}$$
 and $u_n(x)=\left(1+\frac{\sin(n^2x^2)}{n}\right)u(x)$. The sequence $(u_n)_{n\in\mathbb{N}^*}$ converges to u in $L^2(\mathbb{R})$ and goes to infinity in $H^1(\mathbb{R})$.

Definition-Theorem (bounded linear operator). A bounded operator on $\mathcal H$ is a linear map $A:\mathcal H\to\mathcal H$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

Definition-Theorem (bounded linear operator). A bounded operator on $\mathcal H$ is a linear map $A:\mathcal H\to\mathcal H$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

The set $\mathcal{B}(\mathcal{H})$ of the bounded operators on \mathcal{H} is a non-commutative algebra and $\|\cdot\|$ is a norm on $\mathcal{B}(\mathcal{H})$.

Definition-Theorem (bounded linear operator). A bounded operator on $\mathcal H$ is a linear map $A:\mathcal H\to\mathcal H$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

The set $\mathcal{B}(\mathcal{H})$ of the bounded operators on \mathcal{H} is a non-commutative algebra and $\|\cdot\|$ is a norm on $\mathcal{B}(\mathcal{H})$.

Remark. A bounded linear operator is uniquely defined by the values of the sesquilinear form $\mathcal{H} \times \mathcal{H} \ni (u,v) \mapsto \langle u|Av \rangle \in \mathbb{C}$.

Definition-Theorem (bounded linear operator). A bounded operator on $\mathcal H$ is a linear map $A:\mathcal H\to\mathcal H$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

The set $\mathcal{B}(\mathcal{H})$ of the bounded operators on \mathcal{H} is a non-commutative algebra and $\|\cdot\|$ is a norm on $\mathcal{B}(\mathcal{H})$.

Remark. A bounded linear operator is uniquely defined by the values of the sesquilinear form $\mathcal{H} \times \mathcal{H} \ni (u,v) \mapsto \langle u|Av \rangle \in \mathbb{C}$.

Definition-Theorem (adjoint of a bounded linear operator). Let $A \in \mathcal{B}(\mathcal{H})$. The operator $A^* \in \mathcal{B}(\mathcal{H})$ defined by

$$\forall (u, v) \in \mathcal{H} \times \mathcal{H}, \quad \langle u | A^* v \rangle = \langle A u | v \rangle,$$

is called the adjoint of A.

Definition-Theorem (bounded linear operator). A bounded operator on $\mathcal H$ is a linear map $A:\mathcal H\to\mathcal H$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

The set $\mathcal{B}(\mathcal{H})$ of the bounded operators on \mathcal{H} is a non-commutative algebra and $\|\cdot\|$ is a norm on $\mathcal{B}(\mathcal{H})$.

Remark. A bounded linear operator is uniquely defined by the values of the sesquilinear form $\mathcal{H} \times \mathcal{H} \ni (u,v) \mapsto \langle u|Av \rangle \in \mathbb{C}$.

Definition-Theorem (adjoint of a bounded linear operator). Let $A \in \mathcal{B}(\mathcal{H})$. The operator $A^* \in \mathcal{B}(\mathcal{H})$ defined by

$$\forall (u, v) \in \mathcal{H} \times \mathcal{H}, \quad \langle u | A^* v \rangle = \langle A u | v \rangle,$$

is called the adjoint of A. The operator A is called self-adjoint if $A^* = A$.

Definition-Theorem (bounded linear operator). A bounded operator on \mathcal{H} is a linear map $A: \mathcal{H} \to \mathcal{H}$ such that

$$||A|| := \sup_{u \in \mathcal{H} \setminus \{0\}} \frac{||Au||}{||u||} < \infty.$$

The set $\mathcal{B}(\mathcal{H})$ of the bounded operators on \mathcal{H} is a non-commutative algebra and $\|\cdot\|$ is a norm on $\mathcal{B}(\mathcal{H})$.

Remark. A bounded linear operator is uniquely defined by the values of the sesquilinear form $\mathcal{H} \times \mathcal{H} \ni (u,v) \mapsto \langle u|Av \rangle \in \mathbb{C}$.

Definition-Theorem (adjoint of a bounded linear operator). Let $A \in \mathcal{B}(\mathcal{H})$. The operator $A^* \in \mathcal{B}(\mathcal{H})$ defined by

$$\forall (u, v) \in \mathcal{H} \times \mathcal{H}, \quad \langle u | A^* v \rangle = \langle A u | v \rangle,$$

is called the adjoint of A. The operator A is called self-adjoint if $A^* = A$.

Endowed with its norm $\|\cdot\|$ and the * operation, $\mathcal{B}(\mathcal{H})$ is a C*-algebra.

43

(Non necessarily bounded) linear operators on Hilbert spaces

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A.

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A. Note that bounded linear operators are particular linear operators.

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A. Note that bounded linear operators are particular linear operators.

Definition (extensions of operators). Let A_1 and A_2 be operators on \mathcal{H} . A_2 is called an extension of A_1 if $D(A_1) \subset D(A_2)$ and if $\forall u \in D(A_1)$, $A_2u = A_1u$.

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A. Note that bounded linear operators are particular linear operators.

Definition (extensions of operators). Let A_1 and A_2 be operators on \mathcal{H} . A_2 is called an extension of A_1 if $D(A_1) \subset D(A_2)$ and if $\forall u \in D(A_1)$, $A_2u = A_1u$.

Definition (unbounded linear operator). An operator A on \mathcal{H} which does not possess a bounded extension is called an unbounded operator on \mathcal{H} .

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A. Note that bounded linear operators are particular linear operators.

Definition (extensions of operators). Let A_1 and A_2 be operators on \mathcal{H} . A_2 is called an extension of A_1 if $D(A_1) \subset D(A_2)$ and if $\forall u \in D(A_1)$, $A_2u = A_1u$.

Definition (unbounded linear operator). An operator A on \mathcal{H} which does not possess a bounded extension is called an unbounded operator on \mathcal{H} .

Definition (symmetric operator). A linear operator A on $\mathcal H$ with dense domain D(A) is called symmetric if

$$\forall (u, v) \in D(A) \times D(A), \quad \langle Au|v \rangle = \langle u|Av \rangle.$$

Definition (linear operator). A linear operator on \mathcal{H} is a linear map $A:D(A)\to\mathcal{H}$, where D(A) is a subspace of \mathcal{H} called the domain of A. Note that bounded linear operators are particular linear operators.

Definition (extensions of operators). Let A_1 and A_2 be operators on \mathcal{H} . A_2 is called an extension of A_1 if $D(A_1) \subset D(A_2)$ and if $\forall u \in D(A_1)$, $A_2u = A_1u$.

Definition (unbounded linear operator). An operator A on \mathcal{H} which does not possess a bounded extension is called an unbounded operator on \mathcal{H} .

Definition (symmetric operator). A linear operator A on $\mathcal H$ with dense domain D(A) is called symmetric if

$$\forall (u, v) \in D(A) \times D(A), \quad \langle Au|v \rangle = \langle u|Av \rangle.$$

Symmetric operators are not very interesting. Only self-adjoint operators represent physical observables and have nice mathematical properties:

- real spectrum;
- spectral decomposition and functional calculus.

$$D(A^*) = \{ v \in \mathcal{H} \mid \exists w_v \in \mathcal{H} \text{ s.t. } \forall u \in D(A), \ \langle Au | v \rangle = \langle u | w_v \rangle \}.$$

The linear operator A^* on \mathcal{H} , with domain $D(A^*)$, defined by

$$\forall v \in D(A^*), \quad A^*v = w_v,$$

(if w_v exists, it is unique since D(A) is dense) is called the adjoint of A.

$$D(A^*) = \{ v \in \mathcal{H} \mid \exists w_v \in \mathcal{H} \text{ s.t. } \forall u \in D(A), \ \langle Au | v \rangle = \langle u | w_v \rangle \}.$$

The linear operator A^* on \mathcal{H} , with domain $D(A^*)$, defined by

$$\forall v \in D(A^*), \quad A^*v = w_v,$$

(if w_v exists, it is unique since D(A) is dense) is called the adjoint of A. (This definition agrees with the one on Slide 6 for bounded operators.)

$$D(A^*) = \{ v \in \mathcal{H} \mid \exists w_v \in \mathcal{H} \text{ s.t. } \forall u \in D(A), \ \langle Au | v \rangle = \langle u | w_v \rangle \}.$$

The linear operator A^* on \mathcal{H} , with domain $D(A^*)$, defined by

$$\forall v \in D(A^*), \quad A^*v = w_v,$$

(if w_v exists, it is unique since D(A) is dense) is called the adjoint of A. (This definition agrees with the one on Slide 6 for bounded operators.)

Definition (self-adjoint operator). A linear operator A with dense domain is called self-adjoint if $A^* = A$ (that is if A symmetric and $D(A^*) = D(A)$).

$$D(A^*) = \{ v \in \mathcal{H} \mid \exists w_v \in \mathcal{H} \text{ s.t. } \forall u \in D(A), \ \langle Au | v \rangle = \langle u | w_v \rangle \}.$$

The linear operator A^* on \mathcal{H} , with domain $D(A^*)$, defined by

$$\forall v \in D(A^*), \quad A^*v = w_v,$$

(if w_v exists, it is unique since D(A) is dense) is called the adjoint of A. (This definition agrees with the one on Slide 6 for bounded operators.)

Definition (self-adjoint operator). A linear operator A with dense domain is called self-adjoint if $A^* = A$ (that is if A symmetric and $D(A^*) = D(A)$).

Case of bounded operators:

symmetric \Leftrightarrow self-adjoint.

$$D(A^*) = \{ v \in \mathcal{H} \mid \exists w_v \in \mathcal{H} \text{ s.t. } \forall u \in D(A), \ \langle Au | v \rangle = \langle u | w_v \rangle \}.$$

The linear operator A^* on \mathcal{H} , with domain $D(A^*)$, defined by

$$\forall v \in D(A^*), \quad A^*v = w_v,$$

(if w_v exists, it is unique since D(A) is dense) is called the adjoint of A. (This definition agrees with the one on Slide 6 for bounded operators.)

Definition (self-adjoint operator). A linear operator A with dense domain is called self-adjoint if $A^* = A$ (that is if A symmetric and $D(A^*) = D(A)$).

Case of bounded operators:

symmetric \Leftrightarrow self-adjoint.

Case of unbounded operators:

symmetric (easy to check) $\stackrel{\Rightarrow}{\Leftarrow}$ self-adjoint (sometimes difficult to check)

Some unbounded self-adjoint operators arising in quantum mechanics

- \bullet position operator along the j axis:
 - $-\mathcal{H}=L^2(\mathbb{R}^d),$
 - $-D(\widehat{r}_j) = \left\{ u \in L^2(\mathbb{R}^d) \mid r_j u \in L^2(\mathbb{R}^d) \right\}, (\widehat{r}_j \phi)(\mathbf{r}) = r_j \phi(\mathbf{r});$
- \bullet momentum operator along the j axis:
 - $-\mathcal{H}=L^2(\mathbb{R}^d)$,
 - $-D(\widehat{p}_j) = \left\{ u \in L^2(\mathbb{R}^d) \mid \partial_{r_j} u \in L^2(\mathbb{R}^d) \right\}, (\widehat{p}_j \phi)(\mathbf{r}) = -i \partial_{r_j} \phi(\mathbf{r});$
- kinetic energy operator:
 - $-\mathcal{H}=L^2(\mathbb{R}^d),$
 - $-D(T) = H^2(\mathbb{R}^d) := \{ u \in L^2(\mathbb{R}^d) \mid \Delta u \in L^2(\mathbb{R}^d) \}, T = -\frac{1}{2}\Delta;$
- Schrödinger operators in 3D: let $V \in L^2_{\mathrm{unif}}(\mathbb{R}^3,\mathbb{R})$ ($V(\mathbf{r}) = -\frac{Z}{|\mathbf{r}|}$ OK)
 - $-\mathcal{H}=L^2(\mathbb{R}^3),$
 - $-D(H) = H^2(\mathbb{R}^3), H = -\frac{1}{2}\Delta + V.$

Definition-Theorem (spectrum of a linear operator). Let A be a closed¹ linear operator on \mathcal{H} .

• The open set $\rho(A)=\{z\in\mathbb{C}\mid (z-A):D(A)\to\mathcal{H} \text{ invertible}\}$ is called the resolvent set of A.

 $^{^1}$ The operator A is called closed if its graph $\Gamma(A):=\{(u,Au),\,u\in D(A)\}$ is a closed subspace of $\mathcal{H}\times\mathcal{H}$.

Definition-Theorem (spectrum of a linear operator). Let A be a closed¹ linear operator on \mathcal{H} .

• The open set $\rho(A)=\{z\in\mathbb{C}\mid (z-A):D(A)\to\mathcal{H} \text{ invertible}\}$ is called the resolvent set of A. The analytic function

$$\rho(A) \ni z \mapsto R_z(A) := (z - A)^{-1} \in \mathcal{B}(\mathcal{H})$$

is called the resolvent of A. It holds $R_z(A) - R_{z'}(A) = (z'-z)R_z(A)R_{z'}(A)$.

 $^{^{1} \}text{ The operator } A \text{ is called closed if its graph } \Gamma(A) := \{(u,Au), \ u \in D(A)\} \text{ is a closed subspace of } \mathcal{H} \times \mathcal{H}.$

Definition-Theorem (spectrum of a linear operator). Let A be a closed linear operator on \mathcal{H} .

• The open set $\rho(A)=\{z\in\mathbb{C}\mid (z-A):D(A)\to\mathcal{H} \text{ invertible}\}$ is called the resolvent set of A. The analytic function

$$\rho(A) \ni z \mapsto R_z(A) := (z - A)^{-1} \in \mathcal{B}(\mathcal{H})$$

is called the resolvent of A. It holds $R_z(A)-R_{z'}(A)=(z'-z)R_z(A)R_{z'}(A)$.

ullet The closed set $\sigma(A)=\mathbb{C}\setminus
ho(A)$ is called the spectrum of A.

 $^{^{1} \}text{ The operator } A \text{ is called closed if its graph } \Gamma(A) := \{(u,Au), \ u \in D(A)\} \text{ is a closed subspace of } \mathcal{H} \times \mathcal{H}.$

Definition-Theorem (spectrum of a linear operator). Let A be a closed linear operator on \mathcal{H} .

• The open set $\rho(A)=\{z\in\mathbb{C}\mid (z-A):D(A)\to\mathcal{H} \text{ invertible}\}$ is called the resolvent set of A. The analytic function

$$\rho(A) \ni z \mapsto R_z(A) := (z - A)^{-1} \in \mathcal{B}(\mathcal{H})$$

is called the resolvent of A. It holds $R_z(A)-R_{z'}(A)=(z'-z)R_z(A)R_{z'}(A)$.

- ullet The closed set $\sigma(A)=\mathbb{C}\setminus
 ho(A)$ is called the spectrum of A.
- If A is self-adjoint, then $\sigma(A) \subset \mathbb{R}$

 $^{^{1} \}text{ The operator } A \text{ is called closed if its graph } \Gamma(A) := \{(u,Au), \ u \in D(A)\} \text{ is a closed subspace of } \mathcal{H} \times \mathcal{H}.$

Definition-Theorem (spectrum of a linear operator). Let A be a closed¹ linear operator on \mathcal{H} .

• The open set $\rho(A) = \{z \in \mathbb{C} \mid (z - A) : D(A) \to \mathcal{H} \text{ invertible} \}$ is called the resolvent set of A. The analytic function

$$\rho(A) \ni z \mapsto R_z(A) := (z - A)^{-1} \in \mathcal{B}(\mathcal{H})$$

is called the resolvent of A. It holds $R_z(A)-R_{z'}(A)=(z'-z)R_z(A)R_{z'}(A)$.

- The closed set $\sigma(A) = \mathbb{C} \setminus \rho(A)$ is called the spectrum of A.
- If A is self-adjoint, then $\sigma(A) \subset \mathbb{R}$ and it holds $\sigma(A) = \sigma_p(A) \cup \sigma_c(A)$, where $\sigma_p(A)$ and $\sigma_c(A)$ are respectively the point spectrum and the continuous spectrum of A defined as

$$\sigma_{\mathrm{p}}(A) \ = \ \{z \in \mathbb{C} \mid (z-A) \ : \ D(A) \to \mathcal{H} \ \text{non-injective}\} = \{\text{eigenvalues of} \ A\}$$

$$\sigma_{\mathrm{c}}(A) = \overline{\{z \in \mathbb{C} \mid (z - A) : D(A) \to \mathcal{H} \text{ injective but non surjective}\}}.$$

 $^{^{1} \}text{ The operator } A \text{ is called closed if its graph } \Gamma(A) := \{(u,Au), \ u \in D(A)\} \text{ is a closed subspace of } \mathcal{H} \times \mathcal{H}.$

Theorem (RAGE, Ruelle '69, Amrein and Georgescu '73, Enss '78).

Let H be a locally compact self-adjoint operator on $L^2(\mathbb{R}^d)$.

[Ex.: the Hamiltonian of the hydrogen atom satisfies these assumptions.]

Theorem (RAGE, Ruelle '69, Amrein and Georgescu '73, Enss '78).

Let H be a locally compact self-adjoint operator on $L^2(\mathbb{R}^d)$.

[Ex.: the Hamiltonian of the hydrogen atom satisfies these assumptions.]

Let $\mathcal{H}_p = \overline{\operatorname{Span}\left\{\text{eigenvectors of }H\right\}}$ and $\mathcal{H}_c = \mathcal{H}_p^{\perp}$.

[Ex.: for the Hamiltonian of the hydrogen atom, $\dim(\mathcal{H}_p) = \dim(\mathcal{H}_c) = \infty$.]

Theorem (RAGE, Ruelle '69, Amrein and Georgescu '73, Enss '78).

Let H be a locally compact self-adjoint operator on $L^2(\mathbb{R}^d)$.

[Ex.: the Hamiltonian of the hydrogen atom satisfies these assumptions.]

Let $\mathcal{H}_p = \overline{\operatorname{Span}\left\{\text{eigenvectors of }H\right\}}$ and $\mathcal{H}_c = \mathcal{H}_p^{\perp}$.

[Ex.: for the Hamiltonian of the hydrogen atom, $\dim(\mathcal{H}_p) = \dim(\mathcal{H}_c) = \infty$.]

Let χ_{B_R} be the characteristic function of the ball $B_R = \{\mathbf{r} \in \mathbb{R}^d \mid |\mathbf{r}| < R\}$.

Theorem (RAGE, Ruelle '69, Amrein and Georgescu '73, Enss '78).

Let H be a locally compact self-adjoint operator on $L^2(\mathbb{R}^d)$.

[Ex.: the Hamiltonian of the hydrogen atom satisfies these assumptions.]

Let $\mathcal{H}_p = \overline{\operatorname{\mathbf{Span}}\left\{ \operatorname{\mathbf{eigenvectors}} \ \operatorname{\mathbf{of}} \ H
ight\}} \ \operatorname{\mathbf{and}} \ \mathcal{H}_c = \mathcal{H}_p^\perp.$

[Ex.: for the Hamiltonian of the hydrogen atom, $\dim(\mathcal{H}_p) = \dim(\mathcal{H}_c) = \infty$.]

Let χ_{B_R} be the characteristic function of the ball $B_R = \{\mathbf{r} \in \mathbb{R}^d \mid |\mathbf{r}| < R\}$.

Then

$$(\phi_0 \in \mathcal{H}_p) \Leftrightarrow \forall \varepsilon > 0, \ \exists R > 0, \ \forall t \ge 0, \ \left\| (1 - \chi_{B_R}) e^{-itH} \phi_0 \right\|_{L^2}^2 \le \varepsilon;$$

$$(\phi_0 \in \mathcal{H}_c) \Leftrightarrow \forall R > 0, \lim_{T \to +\infty} \frac{1}{T} \int_0^T \|\chi_{B_R} e^{-itH} \phi_0\|_{L^2}^2 dt = 0.$$

Theorem (RAGE, Ruelle '69, Amrein and Georgescu '73, Enss '78).

Let H be a locally compact self-adjoint operator on $L^2(\mathbb{R}^d)$.

[Ex.: the Hamiltonian of the hydrogen atom satisfies these assumptions.]

Let $\mathcal{H}_p = \overline{\operatorname{Span}\left\{ \operatorname{eigenvectors\ of}\ H \right\}}$ and $\mathcal{H}_c = \mathcal{H}_p^{\perp}$.

[Ex.: for the Hamiltonian of the hydrogen atom, $\dim(\mathcal{H}_p) = \dim(\mathcal{H}_c) = \infty$.]

Let χ_{B_R} be the characteristic function of the ball $B_R = \{\mathbf{r} \in \mathbb{R}^d \mid |\mathbf{r}| < R\}$.

Then

$$(\phi_0 \in \mathcal{H}_p) \Leftrightarrow \forall \varepsilon > 0, \ \exists R > 0, \ \forall t \ge 0, \ \left\| (1 - \chi_{B_R}) e^{-itH} \phi_0 \right\|_{L^2}^2 \le \varepsilon;$$

$$(\phi_0 \in \mathcal{H}_c) \Leftrightarrow \forall R > 0, \lim_{T \to +\infty} \frac{1}{T} \int_0^T \|\chi_{B_R} e^{-itH} \phi_0\|_{L^2}^2 dt = 0.$$

 \mathcal{H}_{p} : set of bound states, \mathcal{H}_{c} : set of scattering states

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

Dirac's bra-ket notation:
$$A = \sum_{n \in \mathbb{N}} \lambda_n |e_n\rangle \langle e_n|, \quad \lambda_n \in \mathbb{R}, \quad \langle e_m|e_n\rangle = \delta_{mn}.$$

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

Dirac's bra-ket notation:
$$A=\sum_{n\in\mathbb{N}}\lambda_n|e_n\rangle\langle e_n|,\quad \lambda_n\in\mathbb{R},\quad \langle e_m|e_n\rangle=\delta_{mn}.$$
 Then,

• the operator A is bounded if and only if $||A|| = \sup_n |\lambda_n| < \infty$;

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

Dirac's bra-ket notation:
$$A=\sum_{n\in\mathbb{N}}\lambda_n|e_n
angle\langle e_n|,\quad \lambda_n\in\mathbb{R},\quad \langle e_m|e_n
angle=\delta_{mn}.$$

Then,

- the operator A is bounded if and only if $||A|| = \sup_n |\lambda_n| < \infty$;
- $D(A) = \{|u\rangle = \sum_{n \in \mathbb{N}} u_n |e_n\rangle \mid \sum_{n \in \mathbb{N}} (1 + |\lambda_n|^2) |u_n|^2 < \infty\rangle\};$

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

Dirac's bra-ket notation:
$$A=\sum_{n\in\mathbb{N}}\lambda_n|e_n\rangle\langle e_n|,\quad \lambda_n\in\mathbb{R},\quad \langle e_m|e_n\rangle=\delta_{mn}.$$

Then,

- the operator A is bounded if and only if $||A|| = \sup_n |\lambda_n| < \infty$;
- $D(A) = \{|u\rangle = \sum_{n \in \mathbb{N}} u_n |e_n\rangle \mid \sum_{n \in \mathbb{N}} (1 + |\lambda_n|^2) |u_n|^2 < \infty\rangle\};$
- $\sigma_{\mathrm{p}}(A) = \{\lambda_n\}_{n \in \mathbb{N}}$ and $\sigma_{\mathrm{c}}(A) = \{\text{accumulation points of } \{\lambda_n\}_{n \in \mathbb{N}}\} \setminus \sigma_{\mathrm{p}}(A);$

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

Dirac's bra-ket notation:
$$A = \sum_{n \in \mathbb{N}} \lambda_n |e_n\rangle \langle e_n|, \quad \lambda_n \in \mathbb{R}, \quad \langle e_m|e_n\rangle = \delta_{mn}.$$

Then,

- the operator A is bounded if and only if $||A|| = \sup_n |\lambda_n| < \infty$;
- $D(A) = \{|u\rangle = \sum_{n \in \mathbb{N}} u_n |e_n\rangle \mid \sum_{n \in \mathbb{N}} (1 + |\lambda_n|^2) |u_n|^2 < \infty\rangle\};$
- ullet $\sigma_{\mathrm{p}}(A)=\{\lambda_n\}_{n\in\mathbb{N}}$ and $\sigma_{\mathrm{c}}(A)=\left\{ ext{accumulation points of } \{\lambda_n\}_{n\in\mathbb{N}} \right\} \setminus \sigma_{\mathrm{p}}(A)$;
- $\mathcal{H}_{p} = \mathcal{H}$ and $\mathcal{H}_{c} = \{0\}$ (no scattering states);

Let A be a self-adjoint operator that can be diagonalized in an orthonormal basis $(e_n)_{n\in\mathbb{N}}$ (this is not the case for many useful self-adjoint operators!).

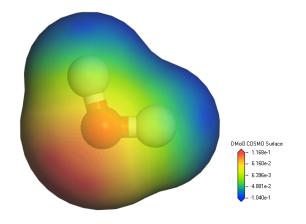
Dirac's bra-ket notation:
$$A = \sum_{n \in \mathbb{N}} \lambda_n |e_n\rangle \langle e_n|, \quad \lambda_n \in \mathbb{R}, \quad \langle e_m|e_n\rangle = \delta_{mn}.$$

Then,

- the operator A is bounded if and only if $||A|| = \sup_n |\lambda_n| < \infty$;
- $D(A) = \{|u\rangle = \sum_{n \in \mathbb{N}} u_n |e_n\rangle \mid \sum_{n \in \mathbb{N}} (1 + |\lambda_n|^2) |u_n|^2 < \infty\rangle\};$
- ullet $\sigma_{\mathrm{p}}(A)=\{\lambda_n\}_{n\in\mathbb{N}}$ and $\sigma_{\mathrm{c}}(A)=\left\{ ext{accumulation points of } \{\lambda_n\}_{n\in\mathbb{N}} \right\} \setminus \sigma_{\mathrm{p}}(A)$;
- $\mathcal{H}_p = \mathcal{H}$ and $\mathcal{H}_c = \{0\}$ (no scattering states);
- ullet functional calculus for diagonalizable self-adjoint operators: for all $f:\mathbb{R} \to \mathbb{C}$, the operator f(A) defined by

$$D(f(A)) = \left\{ |u\rangle = \sum_{n \in \mathbb{N}} u_n |e_n\rangle \mid \sum_{n \in \mathbb{N}} (1 + |f(\lambda_n)|^2) |u_n|^2 < \infty \right\}, \quad f(A) = \sum_{n \in \mathbb{N}} f(\lambda_n) |e_n\rangle \langle e_n|$$

is independent of the choice of the spectral decomposition of A.



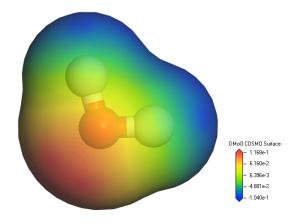
Ex: water molecule H₂O

$$M = 3$$
, $N = 10$, $z_1 = 8$, $z_2 = 1$, $z_3 = 1$

$$v_{ ext{ext}}(\mathbf{r}) = -\sum_{k=1}^{M} \frac{z_k}{|\mathbf{r} - \mathbf{R}_k|}$$

$$\left| \left(-\frac{1}{2} \sum_{i=1}^{N} \Delta_{\mathbf{r}_i} + \sum_{i=1}^{N} v_{\text{ext}}(\mathbf{r}_i) + \sum_{1 \le i < j \le N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right) \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = E \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) \right|$$

$$\forall p \in \mathfrak{S}_N, \quad \Psi(\mathbf{r}_{p(1)}, \cdots, \mathbf{r}_{p(N)}) = \varepsilon(p)\Psi(\mathbf{r}_1, \cdots, \mathbf{r}_N),$$
 (Pauli principle)



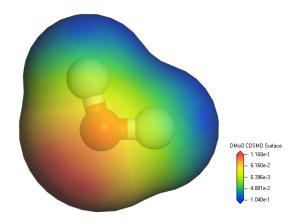
Ex: water molecule H₂O

$$M = 3$$
, $N = 10$, $z_1 = 8$, $z_2 = 1$, $z_3 = 1$

$$v_{ ext{ext}}(\mathbf{r}) = -\sum_{k=1}^{M} \frac{z_k}{|\mathbf{r} - \mathbf{R}_k|}$$

$$\left| \left(-\frac{1}{2} \sum_{i=1}^{N} \Delta_{\mathbf{r}_i} + \sum_{i=1}^{N} v_{\text{ext}}(\mathbf{r}_i) + \sum_{1 \le i < j \le N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right) \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = E \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) \right|$$

$$\Psi \in \mathcal{H}_N = \bigwedge^N \mathcal{H}_1, \qquad \mathcal{H}_1 = L^2(\mathbb{R}^3, \mathbb{C})$$



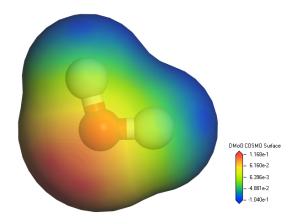
Ex: water molecule H₂O

$$M = 3$$
, $N = 10$, $z_1 = 8$, $z_2 = 1$, $z_3 = 1$

$$v_{ ext{ext}}(\mathbf{r}) = -\sum_{k=1}^{M} \frac{z_k}{|\mathbf{r} - \mathbf{R}_k|}$$

$$\left| \left(-\frac{1}{2} \sum_{i=1}^{N} \Delta_{\mathbf{r}_i} + \sum_{i=1}^{N} v_{\text{ext}}(\mathbf{r}_i) + \sum_{1 \le i < j \le N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right) \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = E \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) \right|$$

$$\Psi \in \mathcal{H}_N = \bigwedge^N \mathcal{H}_1, \qquad \mathcal{H}_1 = L^2(\mathbb{R}^3, \mathbb{C}^2)$$
 with spin



Ex: water molecule H₂O

$$M = 3$$
, $N = 10$, $z_1 = 8$, $z_2 = 1$, $z_3 = 1$

$$v_{\mathrm{ext}}(\mathbf{r}) = -\sum_{k=1}^{M} \frac{z_k}{|\mathbf{r} - \mathbf{R}_k|}$$

$$\left| \left(-\frac{1}{2} \sum_{i=1}^{N} \Delta_{\mathbf{r}_i} + \sum_{i=1}^{N} v_{\text{ext}}(\mathbf{r}_i) + \sum_{1 \le i < j \le N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right) \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = E \Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) \right|$$

$$\Psi \in \mathcal{H}_N = \bigwedge^N \mathcal{H}_1, \qquad \mathcal{H}_1 = L^2(\mathbb{R}^3, \mathbb{C})$$
 Theorem (Kato '51). The operator $H_N := -\frac{1}{2} \sum_{i=1}^N \Delta_{\mathbf{r}_i} + \sum_{i=1}^N v_{\mathrm{ext}}(\mathbf{r}_i) + \sum_{1 \leq i < j \leq N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|}$ with domain $D(H_N) := \mathcal{H}_N \cap H^2(\mathbb{R}^{3N})$ is self-adjoint on \mathcal{H}_N .

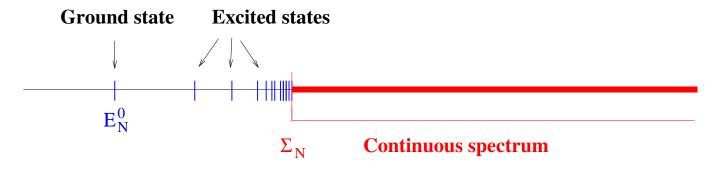
Theorem (spectrum of H_N).

1. HVZ theorem (Hunziger '66, van Winten '60, Zhislin '60)

$$\sigma_{\rm c}(H_N) = [\Sigma_N, +\infty)$$
 with $\Sigma_N = \min \sigma(H_{N-1}) \le 0$ and $\Sigma_N < 0$ iff $N \ge 2$.

2. Bound states of neutral molecules and positive ions (Zhislin '61)

If $N \leq Z := \sum_{k=1}^{M} z_k$, then H_N has an infinite number of bound states.



3. Bound states of negative ions (Yafaev '72)

If $N \geq Z+1$, then H_N has at most a finite number of bound states.

Spectra of Schrödinger operators with confining potentials

$$\mathcal{H}=L^2(\mathbb{R}^d), \qquad V\in C^0(\mathbb{R}^d), \qquad \lim_{|\mathbf{r}|\to +\infty}V(\mathbf{r})=+\infty \ ext{(confining potential)}$$

$$D(H) = \left\{u \in L^2(\mathbb{R}^d) \mid \; -\frac{1}{2}\Delta u + Vu \in L^2(\mathbb{R}^d)\right\}, \quad \forall u \in D(H), \; Hu = -\frac{1}{2}\Delta u + Vu.$$

H is bounded below and its spectrum is purely discrete ($\sigma_{\rm d}(H)=\sigma(H)$, $\sigma_{\rm c}(H)=\emptyset$).

As a consequence, H is diagonalizable in a orthonormal basis: there exist

- a non-decreasing sequence $(E_n)_{n\in\mathbb{N}}$ of real numbers going to $+\infty$;
- ullet an orthonormal basis $(\psi_n)_{n\in\mathbb{N}}$ of \mathcal{H} composed of vectors of D(H), such that

$$\forall n \in \mathbb{N}, \quad H\psi_n = E_n\psi_n.$$

In addition, the ground state eigenvalue E_0 is non-degenerate and the corresponding eigenvector can be chosen positive on \mathbb{R}^d .

Spectra of 3D Schrödinger operators with potentials decaying at infinity

$$V \quad \text{such that} \quad \forall \varepsilon > 0, \ \exists (V_2, V_\infty) \in L^2(\mathbb{R}^3) \times L^\infty(\mathbb{R}^3) \text{ s.t. } V = V_2 + V_\infty \text{ and } \|V_\infty\|_{L^\infty} \leq \varepsilon,$$

$$\mathcal{H} = L^2(\mathbb{R}^3), \qquad D(H) = H^2(\mathbb{R}^3), \qquad \forall u \in D(H), \ Hu = -\frac{1}{2}\Delta u + Vu.$$

The operator H is self-adjoint, bounded below, and $\sigma_{\rm c}(H)=[0,+\infty)$.

Depending on V, the discrete spectrum of H may be

- the empty set;
- a finite number of negative eigenvalues;
- a countable infinite number of negative eigenvalues accumulating at 0 (ex: Ridberg states).

If H has a ground state, then its energy is a non-degenerate eigenvalue and the corresponding eigenvector can be chosen positive on \mathbb{R}^d .

The special case of Kohn-Sham LDA Hamiltonians

$$H_{\rho} = -\frac{1}{2}\Delta + V_{\rho}^{\mathrm{KS}} \quad \text{with} \quad V_{\rho}^{\mathrm{KS}}(\mathbf{r}) = -\sum_{k=1}^{M} \frac{z_{k}}{|\mathbf{r} - \mathbf{R}_{k}|} + \int_{\mathbb{R}^{3}} \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}' + \frac{de_{\mathrm{xc}}^{\mathrm{LDA}}}{d\rho}(\rho(\mathbf{r}))$$

For any $\rho \in L^1(\mathbb{R}^3) \cap L^3(\mathbb{R}^3)$, the KS potential V_{ρ}^{KS} satisfies the assumptions of the previous slide. In particular H_{ρ} is bounded below and $\sigma_{\rm c}(H_{\rho}) = [0, +\infty)$.

Let
$$Z=\sum_{k=1}^M z_k$$
 be the total nuclear charge of the molecular system and $N=\int_{\mathbb{R}^3}
ho$.

- If N < Z (positive ion), H_{ρ} has a countable infinite number of negative eigenvalues accumulating at 0.
- If N=Z (neutral molecular system) and if ρ is a ground state density of the system, then H_{ρ} has at least N non-positive eigenvalues.

Spectra of Hartree-Fock Hamiltonians

$$\begin{aligned} \text{Let } \Phi &= (\phi_1, \cdots, \phi_N) \in (H^1(\mathbb{R}^3))^N \text{ be such that } \int_{\mathbb{R}^3} \phi_i \phi_j = \delta_{ij}, \\ \gamma(\mathbf{r}, \mathbf{r}') &= \sum_{i=1}^N \phi_i(\mathbf{r}) \phi_i(\mathbf{r}'), \qquad \rho_\gamma(\mathbf{r}) = \gamma(\mathbf{r}, \mathbf{r}) = \sum_{i=1}^N |\phi_i(\mathbf{r})|^2. \\ \mathcal{H} &= L^2(\mathbb{R}^3), \quad D(H) = H^2(\mathbb{R}^3), \\ (H\phi)(\mathbf{r}) &= -\frac{1}{2} \Delta \phi(\mathbf{r}) - \sum_{l=1}^M \frac{z_k}{|\mathbf{r} - \mathbf{R}_k|} \phi(\mathbf{r}) + \left(\int_{\mathbb{R}^3} \frac{\rho_\gamma(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \, d\mathbf{r}' \right) \phi(\mathbf{r}) - \int_{\mathbb{R}^3} \frac{\gamma(\mathbf{r}, \mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \, \phi(\mathbf{r}') \, d\mathbf{r}' \end{aligned}$$

Let $Z := \sum_{k=1}^{M} z_k$. The operator H is self-adjoint, bounded below, and we have:

- $\sigma_{\rm ess} = [0, +\infty)$;
- ullet if N < Z (positive ion), H has a countable infinite number of negative eigenvalues accumulating at 0;
- ullet if N=Z (neutral molecular system) and if Φ is a HF minimizer of the system, then H has at least N negative eigenvalues (counting multiplicities).

Spectra of Dirac Hamiltonians

$$\mathcal{H} = L^2(\mathbb{R}^3; \mathbb{C}^4), \qquad D(D_0) = H^1(\mathbb{R}^3; \mathbb{C}^4), \qquad D_0 = c\vec{p} \cdot \vec{\alpha} + mc^2\beta$$

$$p_j = -i\hbar \partial_j, \qquad \alpha_j = \begin{pmatrix} 0 & \sigma_k \\ \sigma_k & 0 \end{pmatrix} \in \mathbb{C}^{4\times 4}, \qquad \beta = \begin{pmatrix} I_2 & 0 \\ 0 & -I_2 \end{pmatrix} \in \mathbb{C}^{4\times 4}$$

$$\sigma_1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \qquad \sigma_2 = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix}, \qquad \sigma_3 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$
 (Pauli matrices)

The free Dirac operator D_0 is self-adjoint and

$$\sigma(D_0) = \sigma_{ac}(D_0) = (-\infty, -mc^2] \cup [mc^2, +\infty).$$

Theorem. Let $\alpha:=\frac{e^2}{4\pi\varepsilon_0\hbar c}\simeq 1/137.036$ be the fine structure constant. Let

$$D_Z = D_0 - \frac{Z}{|\mathbf{r}|}, \qquad Z \in \mathbb{R} \quad \text{(physical cases: } Z = 1, 2, 3, \cdots \text{)}.$$

- if $|Z| < \frac{\sqrt{3}}{2\alpha} \simeq 118.677$, the Dirac operator D_Z is essentially self-adjoint (meaning that there exists a unique domain $D(D_Z)$ containing $C_{\rm c}^{\infty}(\mathbb{R}^3;\mathbb{C}^4)$ for which D_Z is self-adjoint);
- if $|Z| > \frac{\sqrt{3}}{2\alpha} \simeq 118.677$, D_Z has many self-adjoint extensions;
- if $|Z|<\frac{1}{\alpha}\simeq 137.036$, D_Z has a special self-adjoint extension, considered as the physical one. The essential spectrum of this self-adjoint extension is $(-\infty,-mc^2]\cup[mc^2,+\infty)$ and its discrete spectrum consist of the eigenvalues

$$E_{nj} = mc^{2} \left[1 + \left(\frac{Z\alpha}{n - j - \frac{1}{2} + \sqrt{(j + \frac{1}{2})^{2} - Z^{2}\alpha^{2}}} \right)^{2} \right], \quad n \in \mathbb{N}^{*}, \ j = \frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \dots \leq n - \frac{1}{2}.$$

Many-body Dirac-Coulomb Hamiltonian are not understood mathematically.

Theorem (functional calculus for bounded functions). Let $\mathfrak{B}(\mathbb{R},\mathbb{C})$ be the *-algebra of bounded \mathbb{C} -valued Borel functions on \mathbb{R} and let A be a self-adjoint operator on \mathcal{H} . Then there exists a unique map

$$\Phi_A: \mathfrak{B}(\mathbb{R},\mathbb{C})\ni f\mapsto f(A)\in \mathcal{B}(\mathcal{H})$$

satisfies the following properties:

1. Φ_A is a homomorphism of *-algebras:

$$(\alpha f + \beta g)(A) = \alpha f(A) + \beta g(A), \quad (fg)(A) = f(A)g(A), \quad \overline{f}(A) = f(A)^*;$$

- **2.** $||f(A)|| \le \sup_{x \in \mathbb{R}} |f(x)|;$
- 3. if $f_n(x) \to x$ pointwise and $|f_n(x)| \le |x|$ for all n and all $x \in \mathbb{R}$, then $\forall u \in D(A), \quad f_n(A)u \to Au$ in \mathcal{H} ;
- **4.** if $f_n(x) \to f(x)$ pointwise and $\sup_n \sup_{x \in \mathbb{R}} |f_n(x)| < \infty$, then $\forall u \in \mathcal{H}, \quad f_n(A)u \to f(A)u \text{ in } \mathcal{H};$

In addition, if $u \in \mathcal{H}$ is such that $Au = \lambda u$, then $f(A)u = f(\lambda)u$.

Theorem (spectral projections and functional calculus - general case -).

Let A be a self-adjoint operator on \mathcal{H} .

- For all $\lambda \in \mathbb{R}$, the bounded operator $P_{\lambda}^{A} := \mathbb{1}_{]-\infty,\lambda]}(A)$, where $\mathbb{1}_{]-\infty,\lambda]}(\cdot)$ is the characteristic function of $]-\infty,\lambda]$, is an orthogonal projection.
- Spectral decomposition of A: for all $u \in D(A)$ and $v \in \mathcal{H}$, it holds

$$\langle v|Au \rangle = \int_{\mathbb{R}} \lambda \, \underline{d\langle v|P_{\lambda}^{A}u \rangle}, \quad \text{which we denote by} \quad A = \int_{\mathbb{R}} \lambda \, dP_{\lambda}^{A}.$$
Bounded complex measure on \mathbb{R}

• Functional calculus: let f be a (not necessarily bounded) \mathbb{C} -valued Borel function on \mathbb{R} . The operator f(A) can be defined by

$$D(f(A)) := \left\{ u \in \mathcal{H} \mid \int_{\mathbb{R}} |f(\lambda)|^2 \underline{d\langle u| P_{\lambda}^A u \rangle} < \infty \right\}$$

Bounded positive measure on \mathbb{R}

and

$$\forall (u,v) \in D(f(A)) \times \mathcal{H}, \ \langle v|f(A)u \rangle := \int_{\mathbb{R}} f(\lambda) \, d\langle v|P_{\lambda}^{A}u \rangle.$$