# basics - math

**basics**

**Feb 19, 2026**

# CONTENTS

This material is part of the **basics-books project**. It is also available as a .pdf document.

**Contents.**

- **Linear Algebra.** ..., matrix factorization,...; basics of lots of numerical methods; **Vector and Tensor Algebra** provides the mathematical framework for manipulating vectors and tensors, the mathematical objects meant to represent **absolute quantities** - invariant under coordinate transformations - fundamental in geometry and physics.

- **Multivariable Calculus** provides the tools for working with continuous functions of many variables; **Differential Geometry** studies smooth curves, surfaces, volumes (and manifolds in general), within the framework of a Riemann structure, providing a metric to define distances, angles, curvatures and other geometric properties; **Vector and Tensor Calculus** studies vector and tensor fields defined on a manifold, along with their differentiation and integration, and thus being essential for a mature approach to geometry and physics;

- **Functional analysis** and **Complex Calculus** (complex analysis, transforms,...)

- ODEs

- PDEs

- Calculus of Variations: theoretical background; Lagrange multiplier method for constraints; sensitivty; gradient-based methods,...

- Optimization

# Part I

# Linear Algebra

# ONE

# INTRODUCTION TO LINEAR ALGEBRA

- Matrices:

    - properties

    - representation of linear transformations

    - factorizations

    - linear systems

todo:

- existence of solution of equilibria of LTI and their stability, in *LTI:Equilibria* and in *LTI:Stability*

# MATRICES

$\mathbf{A} \in \mathbb{K}^{m,n}$ with usually $\mathbb{K}^{m,n} = \mathbb{R}^{m,n}$ or $\mathbb{C}^{m,n}$

**Hermitian matrix.** The Hermitian matrix $\mathbf{A}^*$ of a matrix $\mathbf{A}$ is the transpose and complex conjugate matrix (if $\mathbb{K} = \mathbb{C}$),

$$[\mathbf{A}^*]_{ij} = A_{ji}^* \,,$$

with the notation of $a^*$ for the complex conjugate of a numerical quantity.

## 2.1 Subspaces

### 2.1.1 Range, Image

$$R(\mathbf{A}) = \{\mathbf{y} \in \mathbb{K}^m \mid \exists \mathbf{x} \in \mathbb{K}^m \,, \text{ s.t. } \mathbf{A}\mathbf{x} = \mathbf{y}\}$$

The range of a matrix $\mathbf{A}$ is the linear space built on the columns of $\mathbf{A}$, since the operation $\mathbf{A}\mathbf{x}$ represents nothing but a linear combination of the columns of matrix.

### 2.1.2 Null, Kernel

$$K(\mathbf{A}) = \{\mathbf{x} \in \mathbb{K}^n \mid \mathbf{A}\mathbf{x} = \mathbf{0}\}$$

## 2.2 Theorem

### 2.2.1 Orthogonality of $R(\mathbf{A})$ and $K(\mathbf{A}^*)$

The following holds,

$$R(\mathbf{A}) \perp K(\mathbf{A}^*) \,,$$

meaning that $\forall \mathbf{u} \in R(\mathbf{A})$ and $\forall \mathbf{v} \in K(\mathbf{A}^*)$, $\mathbf{u}^*\mathbf{v} = 0$.

**Proof.**

$$\mathbf{u} = \mathbf{A}\mathbf{x}$$
$$\mathbf{0} = \mathbf{A}^*\mathbf{v}$$

and thus, premultiplication by $\mathbf{x}^*$ of the second relation gives

$$0 = \mathbf{x}^*\mathbf{0} = \underbrace{\mathbf{x}^*\mathbf{A}^*}_{=(\mathbf{A}\mathbf{x})^*=\mathbf{u}^*}\mathbf{v} = \mathbf{u}^*\mathbf{v}\,.$$

This theorem becomes quite useful, e.g. for constrained linear systems and projections… (e.g. N-S, or other constrained linear systems…)

**todo** add links

# MATRIX FACTORIZATIONS

- **Singular Value Decomposition (SVD)**

- **Spectral decomposition** Eigenvalues, eigenvectors; Jordan canonical formula…

## 3.1 Gershgorin circle theorem

Let $\mathbf{A}$ a complex $n, n$ matrix. Let $R_i = \sum_{j \neq i} |a_{ij}|$, and $D(a_{ii}, R) \subseteq \mathbb{C}$ the closed disc centered at $a_{ii}$ with radius $R_i$ in the complex plane.

Every eigenvalue of $\mathbf{A}$ lies within at least one of the Gershgorin discs $D(a_{ii}, R_i)$.

### Proof

Let $\lambda$ an eigenvalue of $\mathbf{A}$. Thus, it exsits a vector $\mathbf{v}$ s.t. $\mathbf{Av} = \lambda \mathbf{v}$, or

$$\sum_j a_{ij} v_j = \lambda v_i \ .$$

Moving $a_{ii}$ from the LHS to the RHS, it follows

$$\sum_{j \neq i} a_{ij} v_j = (\lambda - a_{ii}) v_i \ .$$

Selecting $i$ s.t. $|v_i| \geq |v_j|$, $\forall j$ it follows that

$$|\lambda - a_{ii}| = \left| \sum_{j \neq i} a_{ij} \frac{v_j}{v_i} \right| \leq \sum_{j \neq i} \left| a_{ij} \frac{v_j}{v_i} \right| \leq \sum_{j \neq i} |a_{ij}| \leq R_i \ .$$

## 3.2 Spectral radius

Spectral radius of a matrix $\mathbf{A}$ is defined as the maximum of the absolute value of its eigenvalues

$$\rho(\mathbf{A}) = \max |\lambda(\mathbf{A})| \ .$$

- **QR**
- **LU**
- **Schur**

- **Cholesky** Symmetric positive definite matrices have Choleski decomposition,

$$\mathbf{M} = \mathbf{L}\mathbf{L}^*,$$

with $\mathbf{L}$ lower triangular matrix. And thus quite easy to "invert", for solving linear systems.

## 3.3 Singular Value Decomposition

Singular value decomposition of a matrix $\mathbf{A} \in \mathbb{C}^{m,n}$

$$\mathbf{A}_{(m,n)} = (\text{SVD}) = \mathbf{U}_{(m,m)}\mathbf{S}_{(m,n)}\mathbf{V}^*_{(n,n)}$$

with $\mathbf{U}^*\mathbf{U} = \mathbf{I}_{(m,m)}$, and $\mathbf{V}^*\mathbf{V} = \mathbf{I}_{(n,n)}$, $\mathbf{S} = \text{diag}\{\sigma_i\}$, $\sigma_i \geq 0$.

Exploiting the definition of matrix product, the SVD of matrix A can be written as

$$\mathbf{A} = \sum_{i=\min(m,n)} \sigma_i \mathbf{u}_i \mathbf{v}_i^*,$$

see also economic decomposition below.

### 3.3.1 Properties

**Relation with *range* and *kernel* the matrix.**

$$R(\mathbf{A}) = \{\mathbf{u}_i | \sigma_i > 0\}$$
$$K(\mathbf{A}) = \{\mathbf{v}_i | \sigma_i = 0\}$$
$$R(\mathbf{A}^*) = \{\mathbf{v}_i | \sigma_i > 0\}$$
$$K(\mathbf{A}^*) = \{\mathbf{u}_i | \sigma_i = 0\}$$

It's immediate to prove $R(\mathbf{A}) \perp K(\mathbf{A}^*)$.

**Full or economic decomposition.** In general the $\mathbf{S}$ of the full decompsition in not square.

$$\mathbf{A} = (\text{SVD}) = \mathbf{U}_{(m,m)}\mathbf{S}_{(m,n)}\mathbf{V}^*_{(n,n)} = \mathbf{U}^e_{(m,k)}\mathbf{S}^e_{(k,k)}\mathbf{V}^{e*}_{(k,n)},$$

with $k = \min(m,n)$.

### 3.3.2 Applications

#### Solution of under-determined linear systems

#### Norms and optimization

If $L^2$-norm is used for vector norms, see *Example 3.3.1*

$$\text{Find } \max_{||\mathbf{x}||=1} ||\mathbf{A}\mathbf{x}||^2$$

$$\mathbf{x}^*\mathbf{A}^*\mathbf{A}\mathbf{x} = \mathbf{x}^*\mathbf{V}\mathbf{S}^* \underbrace{\mathbf{U}^*\mathbf{U}}_{\mathbf{I}} \mathbf{S} \underbrace{\mathbf{V}^*\mathbf{x}}_{\mathbf{z}}$$

and $\mathbf{z} = \mathbf{V}^*\mathbf{x}$ is unitary as well (since $1 = \mathbf{x}^*\mathbf{x} = \mathbf{z}^*\mathbf{V}^*\mathbf{V}\mathbf{z} = \mathbf{z}^*\mathbf{z}$).

After defining $\mathbf{S}_2 := \mathbf{S}^*\mathbf{S}$, the problem thus becomes

$$\text{Find } \max_{||\mathbf{z}||=1} \mathbf{z}^*\mathbf{S}_2\mathbf{z}$$

Manipulating the objective function as $\sum_i \sigma_i^2|z_i|^2$, the constraint optimization problem has global maximum $\sigma_1^2$ (sorted singular values from the largest to the smalles) when $\mathbf{z}_1 = (1, 0, 0, \dots, 0)^T$. Going back to the original variable, optimal condition

- is achieved for $\mathbf{x}_1 = \mathbf{v}_1$;

- has value $\max_{||\mathbf{x}||=1} ||\mathbf{A}\mathbf{x}|| = \sigma_1^2$

- and the response of the system is $\mathbf{y}_1 = \sigma_1\mathbf{u}_1$ as

$$\mathbf{y}_1 := \mathbf{A}\mathbf{x}_1 = \mathbf{U}\mathbf{S}\mathbf{V}^*\mathbf{v}_1 = \sum_k \left(\sigma_k\mathbf{u}_k\mathbf{v}_k^*\right)\mathbf{v}_1 = \sigma_1\mathbf{u}_1 \ .$$

---

**Example 3.3.1 (Other norms - variations of the $L^2$-norm)**

This kind of problem may results as the discrete counterpart of a continuous problem, as an example from *finite element methods*, where $\mathbf{x}$, $\mathbf{y}$ contain the coefficients of the basis functions. In this case, the discrete counterpart of the continous norm-measure of the continuous fields may involve a "mass matrix" (symmetric, definite positive - and thus with Choloeski factorization…),

$$\int_{\Omega_x} |x(\vec{r})|^2 d\vec{r} \simeq \mathbf{x}^*\mathbf{M}_x\mathbf{x}$$

$$\int_{\Omega_y} |y(\vec{r})|^2 d\vec{r} \simeq \mathbf{y}^*\mathbf{M}_y\mathbf{y}$$

Continuous and discrete optimization problems are

$$\text{Find } \max_{|x(\vec{r})|_{L^2(\Omega_x)}=1} |y|_{L^2(\Omega_y)}^2$$

$$\text{Find } \max_{\mathbf{x}^*\mathbf{M}_x\mathbf{x}=1} \mathbf{x}^*\mathbf{A}^*\mathbf{M}_y\mathbf{A}\mathbf{x}$$

with the relation $\mathbf{y} = \mathbf{A}\mathbf{x}$ between the discrete input and output.

This problem can be easily (and efficiently?) recast to the standard form of the problem, using *Cholesky decomposition* of matrix $\mathbf{M}_x = \mathbf{L}_x\mathbf{L}_x^*$, with the definition of the variable $\mathbf{z} = \mathbf{L}_x^*\mathbf{x}$

$$\text{Find } \max_{||\mathbf{z}||=1} \mathbf{z}^*\mathbf{L}_x^{-1}\mathbf{A}^*\mathbf{L}_y\mathbf{L}_y^*\mathbf{A}\mathbf{L}_x^{-*}\mathbf{z}$$

This problem can be efficiently solved with iterative algorithms to compute the SVD of the matrix $\widetilde{\mathbf{A}} := \mathbf{L}_y^*\mathbf{A}\mathbf{L}_x^{-*}$, that doesn't need the expensive full inversion of a matrix but only its action on a vector (instead of evaluating the inverse, a linear system - here triangular! Easier to solve - can be efficiently solved). Algorithms like **Arnoldi algorithm** evaluates the largest eigenvalues or singular values(if no options to set other goals) and the corresponding eigenvectors and singular vectors, alternating **power iterations** and **orthogonalization**. Here power iteration to evaluate the action of the matrix $\widetilde{\mathbf{A}} = \mathbf{L}_y^*\mathbf{A}\mathbf{L}_x^{-*}$ on a generic vector $\mathbf{z}$ is made of the following steps:

1. solution of the linear system $\mathbf{L}_x^*\mathbf{a} = \mathbf{z} \to \mathbf{a} = \dots$

2. matrix-vector multiplication $\mathbf{b} = \mathbf{A}\mathbf{a}$

3. matrix-vector multiplication $\mathbf{c} = \mathbf{L}_y^*\mathbf{b}$

Once the SVD is solved, with $\mathbf{z}_1 = \mathbf{v}_1$

$$\mathbf{x}_1 = \mathbf{L}_x^{-*}\mathbf{v}_1$$

$$\mathbf{y}_1 = \mathbf{A}\mathbf{x}_1 = \sigma_1\mathbf{L}_y^{-*}\mathbf{u}_1$$

---

# FOUR

# LINEAR SYSTEMS

The linear system

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

with $\mathbf{A} \in \mathbb{R}^{m,n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$ has solution if there exists (at least) a vector $\tilde{\mathbf{x}}$ whose product with $\mathbf{A}$ gives $\mathbf{b}$.

**Condition for the existence of a solution.** A solution exists if $\mathbf{b}$ belongs to the *range of* $\mathbf{A}$.

**Uniqueness of a solution.** If a solution $\tilde{\mathbf{x}}$ exists, it's unique if the *kernel of* $\mathbf{A}$ is empty, $K(\mathbf{A}) = \emptyset$. If the kernel is not empty,

$$\mathbf{b} = \mathbf{A}(\tilde{\mathbf{x}} + \mathbf{u}) = \mathbf{A}\tilde{\mathbf{x}} + \underbrace{\mathbf{A}\mathbf{u}}_{=\mathbf{0}} \,,$$

for $\forall \mathbf{u} \in K(\mathbf{A})$, and thus an infinite number of solutions exists. Given a vector basis of the kernel $\mathbf{K}(\mathbf{A})$, where $\dim\left(K(\mathbf{A})\right) = n_K$, $\{\mathbf{u}_1, \ldots, \mathbf{u}_{n_K}\}$, the general solution has $n_K$ "degrees of arbitrariness", since the general solution of the problem is

$$\tilde{\mathbf{x}} + \sum_{i=1}^{n_K} a_i \mathbf{u}_i = \tilde{\mathbf{x}} + \mathbf{U}\mathbf{a} \,.$$

**todo** treat under-, det-, over-determined lin sys

# FIVE

# SPECTRAL DECOMPOSITION

- **Introduction.** What's spectral (or eigenvalue) decomposition? Arbitrariness. When is it possible? Special matrices. Properties and theorems (Cayley-Hamilton)

$$\mathbf{A}\mathbf{u}_i = s_i \mathbf{u}_i$$

- **Generalized spectral decompositions.** Generalized spectral decomposition of first order systems; generalized spectral deocmposition of second order systems[1]

$$\mathbf{A}\mathbf{u}_i = s_i \mathbf{B}\mathbf{u}_i$$

$$\left[ s_i^2 \mathbf{M} + s_i \mathbf{C} + \mathbf{K} \right] \mathbf{u}_i = \mathbf{0} \, .$$

- **Sensitivity of spectral decomposition.** …

## 5.1 Spectral decomposition of symmetric matrices

In this section the spectral decomposition of symmetric matrices (or Hermitian, if defined on the complex field), $\mathbf{A} = \mathbf{A}^*$, is discussed.

**Theorem 5.1.1 (Spectral theorem of a symmetric matrix)**

A symmetric matrix has a spectral decomposition, with a complete set of unit orthogonal vectors and real eigenvalues,

$$\mathbf{A} = \mathbf{R}\mathbf{S}\mathbf{R}^* \, ,$$

with $\mathbf{R} = [\, \mathbf{r}_1 \mid \, ... \mid \mathbf{r}_n \,]$, and $\mathbf{r}_i^* \mathbf{r}_k = \delta_{ik}$, and $\mathbf{S} = \mathrm{diag}\{s_i\}$, $s_i \in \mathbb{R}$.

---

[1] It's likely that mainly underdamped systems with simultaneously diagonalizable matrices will be treated, as this kind of systems often arises in structural dynamics of elastic structures with small structural damping.

**Proof**

**todo** *By construction…*

$$\mathbf{A} = \mathbf{RSR}^* =$$

$$= \begin{bmatrix} \mathbf{r}_1 & \dots & \mathbf{r}_n \end{bmatrix} \operatorname{diag}\{s_i\} \begin{bmatrix} \mathbf{r}_1^* \\ \dots \\ \mathbf{r}_n^* \end{bmatrix} =$$

$$= \begin{bmatrix} s_1\mathbf{r}_1 & \dots & s_n\mathbf{r}_n \end{bmatrix} \begin{bmatrix} \mathbf{r}_1^* \\ \dots \\ \mathbf{r}_n^* \end{bmatrix} =$$

$$= s_1\mathbf{r}_1\mathbf{r}_1^* + \dots + s_n\mathbf{r}_n\mathbf{r}_n^* \ .$$

**Properties.**

- Left and right eigenvectors are the same. This immediately follows from the symmetry of the matrix and the definition of the left and right eigenvectors

$$\mathbf{Ar}_k = s_k\mathbf{r}_k$$
$$\mathbf{l}_j^*\mathbf{A} = s_j^*\mathbf{l}_i^* \quad \rightarrow \quad \mathbf{Al}_j = s_j\mathbf{l}_j \ ,$$

- Eigenvalues are real, as

$$0 = \mathbf{l}_k^* \left(\mathbf{Ar}_k\right) - \left(\mathbf{l}_k^*\mathbf{A}\right)\mathbf{r}_k = \left(s_k - s_k^*\right)\mathbf{l}_k^*\mathbf{r}_k = 2\operatorname{im}\{s_k\} \underbrace{\mathbf{l}_k^*\mathbf{r}_k}_{=|\mathbf{r}_k|^2} \ ,$$

and thus $\operatorname{im}\{s_k\} = 0$.

- Eigenvectors with different eigenvalues are orthogonal, as

$$0 = \mathbf{l}_j^* \left(\mathbf{Ar}_k\right) - \left(\mathbf{l}_j^*\mathbf{A}\right)\mathbf{r}_k = \left(s_k - s_j^*\right)\mathbf{l}_j^*\mathbf{r}_k \ ,$$

and thus $\mathbf{l}_j^*\mathbf{r}_k = 0$ for $j \neq k$, if $s_j \neq s_k$. Left and right eigenvectors are orthogonal w.r.t. matrix $\mathbf{A}$ as well.

## 5.2 Sensitivity of spectral decomposition

Matrices involved in a eigenvalue problem can be function of some parameters $p$. Sensitivity analysis evaluates first-order changes[1] of eigenvalues and eigenvectors following a increment of the parameter $p = \overline{p} + \Delta p$,

$$s_i(p) = s_i(\overline{p}) + \Delta p\, s_{i/p}(\overline{p}) + o(\Delta p)$$
$$\mathbf{u}_i(p) = \mathbf{u}_i(\overline{p}) + \Delta p\, \mathbf{u}_{i/p}(\overline{p}) + o(\Delta p) \ .$$

Here the terms $s_{i/p}(\overline{p})$ and $\mathbf{u}_{i/p}(\overline{p})$ are the **sensitivity** of the $i^{th}$ eigenvalue and eigenvector respectively w.r.t. to an increment $\Delta p$ of the variable $p$, starting from the reference configuration determined by the value $\overline{p}$ of the parameter.

**Rank of sensitivity.** Sensitivity of eigenvalue to a scalar parameter is a scalar quantity. Sensitivity of eivenvector to a scalar parameter is a vector quantity. If the parameter $\mathbf{p}$ is a vector quantity, sensitivity of eigenvalue w.r.t. $\mathbf{p}$ is a vector quantity, and sensitivity of eigenvector w.r.t. $\mathbf{p}$ is a matrix (or tensor, sometimes…) quantity.

---

[1] First order changes should be meant as derivatives w.r.t. the parameter $p$, evaluated for the reference value of $\overline{p}$ that appear in a Taylor series of the quantity of interest, $s(p_i) = s(\overline{p}_i) + \Delta p_k\, \partial_{p_k} s(\overline{p}_i) + o(\Delta p_i)$.

## 5.2.1 Generalized eigenvalue problem (first order)

Sensitivity of eigenvalue and right eigenvector to a parameter $p$ of matrices $\mathbf{A}(p)$, $\mathbf{B}(p)$ read

$$s_{i/p} = \mathbf{w}_i^* \left( \mathbf{A}_{/p} - s_i \mathbf{B}_{/p} \right) \mathbf{u}_i \tag{5.1}$$

$$\mathbf{u}_{i/p} = \sum_{a \notin \{i\}} \mathbf{w}_a^* \left( -\mathbf{A}_{/p} + s_i \mathbf{B} \right) \mathbf{u}_i \frac{1}{s_a - s_i} \mathbf{u}_a \ ,$$

having exploited normalization condition $\mathbf{w}_i^* \mathbf{B} \mathbf{u}_i = 1$, being $\mathbf{u}_i$ and $\mathbf{w}_i$ the[2] right and left eigenvectors associated with eigenvalue $s_i$.

### Right and left eigenvalue problem

Right and left eigenvalue problems are defined as

$$\begin{aligned} \mathbf{A}\mathbf{u}_i &= s_i \mathbf{B}\mathbf{u}_i \\ \mathbf{v}_j^* \mathbf{A} &= s_j \mathbf{v}_j^* \mathbf{B} \end{aligned} \tag{5.2}$$

**Property 5.2.1 (Eigenvalues of the right and left spectral problems)**

Left and right eigenvalue problems have the same set of eigenvalues.

### Proof

**todo**

**Property 5.2.2 (Orthogonality conditions of right and left spectral problems)**

Left and right eigenvectors associated with different eigenvalues are orthogonal w.r.t. the matrices of the system, i.e.

$$\begin{aligned} \mathbf{v}_j^* \mathbf{B} \mathbf{u}_i &= b^{(i)} \delta_{ij} \\ \mathbf{v}_j^* \mathbf{A} \mathbf{u}_i &= a^{(i)} \delta_{ij} \ , \end{aligned}$$

with $a^{(i)} = s_i b^{(i)}$. Parameters $a^{(i)}$, $b^{(i)}$ are not uniquely determined, and one of them (or their combination) can be used in a **normalization condition** removing arbitrariness of eigenvectors to a multiplicative factor.

### Proof

Starting from the left and right eigenvalue problems (5.2) for two different eigenvalues $s_i \neq s_j$, and multiplying (on the left) the right eigenvalue problem by $\mathbf{v}_j^*$ and multiplying (on the right) the left eigenvalue problem by $\mathbf{u}_i$

$$\begin{aligned} \mathbf{v}_j^* \mathbf{A} \mathbf{u}_i &= s_i \mathbf{v}_j^* \mathbf{B} \mathbf{u}_i \\ \mathbf{v}_j^* \mathbf{A} \mathbf{u}_i &= s_j \mathbf{v}_j^* \mathbf{B} \mathbf{u}_i \end{aligned} \tag{5.3}$$

---

[2] Eigenvalues with algebraic and geometric multiplicity larger than one, $m_i^{(a)} = m_i^{(g)} > 1$, have the associated eigenvectors that spans a subspace of dimension $m_i$. In this subspace, it's always possible to define a set of orthogonal vectors. **todo** write it better; refer to introduction to spectral decomposition... **todo** treat sensitivity to parameters of eigenvalues and eigenvectors with multiplicity larger than one...

and subracting the two equations

$$0 = (s_i - s_j)\mathbf{v}_j^*\mathbf{B}\mathbf{u}_i$$

$$\rightarrow \qquad 0 = \mathbf{v}_j^*\mathbf{B}\mathbf{u}_i \quad \text{if } s_i \neq s_j , \tag{5.4}$$

as the left-hand-side terms are the same, and $s_i \neq s_j$. As (5.4) holds, from (5.3) it also follows that

$$\rightarrow \qquad 0 = \mathbf{v}_j^*\mathbf{A}\mathbf{u}_i \quad \text{if } s_i \neq s_j . \tag{5.5}$$

These relations don't hold for left and right eigenvectors associated with the same eigenvalues, and thus these conditions can be used as normalization conditions.

## Sensitivity of eigenvalue

The derivative w.r.t. parameter $p$ of the right eigenvalue problem

$$\mathbf{A}(p)\mathbf{u}_i(p) = s_i(p)\mathbf{B}(p)\mathbf{u}_i(p) ,$$

reads

$$0 = (\mathbf{A}_{/p} - s_{i/p}\mathbf{B} - s_i\mathbf{B}_{/p})\mathbf{u}_i + (\mathbf{A} - s_i\mathbf{B})\mathbf{u}_{i/p} .$$

Multiplying on the left by the left eigenvector $\mathbf{w}_i$, and recalling that the left eigenvalue problem reads $0 = \mathbf{w}_i^* (\mathbf{A} - s_i\mathbf{B})$, the last term is identically zero and

$$0 = \mathbf{w}_i^*(\mathbf{A}_{/p} - s_{i/p}\mathbf{B} - s_i\mathbf{B}_{/p})\mathbf{u}_i ,$$

so that it's possible to find the expression of the eigenvalue sensitivity as

$$s_{i/p} = \frac{\mathbf{w}_i\left(\mathbf{A}_{/p} - s_i\mathbf{B}_{/p}\right)\mathbf{u}_i}{\mathbf{w}_i^*\mathbf{B}\mathbf{u}_i} = \mathbf{w}_i\left(\mathbf{A}_{/p} - s_i\mathbf{B}_{/p}\right)\mathbf{u}_i , \tag{5.6}$$

where the last step exploits the normalization condition $\mathbf{w}_i^*\mathbf{B}\mathbf{u}_i = b^{(i)} = 1$.

## Sensitivity of eigenvector

The sensitivity of the eigenvector $\mathbf{u}_{i/p}$ is the solution of the linear system that can be obtained from the derivative of the eigenvalue problem,

$$(\mathbf{A} - s_i\mathbf{B})\mathbf{u}_{i/p} = \underbrace{- (\mathbf{A}_{/p} - s_{i/p}\mathbf{B} - s_i\mathbf{B}_{/p})\mathbf{u}_i}_{=:\mathbf{b}_i} ,$$

once the sensitivity of the eigenvector $s_{i/p}$ is known from expression (5.6).

---

**Property 5.2.3 (Linear system is singular, $\mathbf{K}(A) = \{\mathbf{u}_i\}$)**

Linear system is singular as $s_i$ is an eigenvalue of the problem, $|\mathbf{A} - s_i\mathbf{B}| = 0$, and the kernel of the matrix is the space spanned by the right eigenvectors associated with $s_i$, as

$$(\mathbf{A} - s_i\mathbf{B})\,\mathbf{u}_i = \mathbf{0} .$$

---

As the matrix of the linear system is singular, the linear system has **no solution if** the right-hand side is not in the range of the matrix of the system, i.e. $\mathbf{b}_i \notin \mathrm{R}(\mathbf{A} - s_i\mathbf{B})$.

Fortunately, the RHS belongs to the range of the matrix, $\mathbf{b}_i \in \mathrm{R}(\mathbf{A} - s_i\mathbf{B})$, as it's proved in the box below.

---

**Property 5.2.4 ( $\mathbf{b}_i \in \mathbf{R}(\mathbf{A} - s_i\mathbf{B})$)**

If $\mathbf{b}_i$ belongs to $\mathrm{R}(\mathbf{A} - s_i\mathbf{B})$, it's orthogonal to $\mathrm{K}(\mathbf{A}^* - s_i^*\mathbf{B}^*)$ and viceversa, by theorem … **todo** add link

Kernel of $\mathbf{A}^* - s_i^*\mathbf{B}^*$ is spanned by the left eigenvectors $\mathbf{v}_i$. Thus, using expression of the eigenvalue sensitivity in evaluating the product

$$\mathbf{v}_i^*\mathbf{b}_i = -\mathbf{v}_i^* \left( \mathbf{A}_{/p} - s_{i/p}\mathbf{B} - s_i\mathbf{B}_{/p} \right)\mathbf{u}_i \right) =$$
$$= -\mathbf{v}_i^* \left( \mathbf{A}_{/p} - s_i\mathbf{B}_{/p} \right)\mathbf{u}_i \right) + \mathbf{v}_i^* \left( \mathbf{A}_{/p} - s_i\mathbf{B}_{/p} \right) \mathbf{u}_i \underbrace{\mathbf{v}_i^*\mathbf{B}\mathbf{u}_i}_{=1} =$$
$$= 0 \ ,$$

it's readily proved that $\mathbf{v}_i^*\mathbf{b}_i = 0$. **todo** treat the case where eigenvalue multiplicity is larger than 1.

---

As the singular linear system has at least a solution, it has infinite solution as the kernel of the matrix is not empty. If $\widetilde{\mathbf{u}_{/p}}$ is a solution of the system, adding a linear combination of the vectors of $\mathrm{K}(\mathbf{A} - s_i\mathbf{B}) = \{\mathbf{u}_i\}$ produces a solution as well,

$$\widetilde{\mathbf{u}_{i/p}} + \mathbf{U}_i\beta \ ,$$

as $(\mathbf{A} - s_i\mathbf{B})\mathbf{U}_i = \mathbf{0}$, being $\mathbf{U}_i$ the matrix whose columns are the right eigenvectors $\mathbf{u}_i$ associated with eigenvalue $s_i$, spanning the kernel of the matrix.

In order to remove this arbitrariness, it's possible to introduce the orthogonality condition $\mathbf{V}_i^*\mathbf{B}\widetilde{\mathbf{u}_{i/p}} = \mathbf{0}$, setting the coefficients of the linear combination of the vectors of the kernel $\beta = \mathbf{0}$.

---

**Property 5.2.5 (Augmented linear system)**

$$\begin{bmatrix} \mathbf{A} - s_i\mathbf{B} & \mathbf{U}_i \\ \mathbf{V}_i^*\mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \widetilde{\mathbf{u}_{i/p}} \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{b}_i \\ \mathbf{0} \end{bmatrix}$$

If everything goes right, the coefficients $\beta$ must be zero.

---

**Property 5.2.6 (Decomposition of the soluton using the modal basis)**

Writing the solution as a linear combination of the right eigenvectors,

$$\widetilde{\mathbf{u}_{i/p}} = \mathbf{U}_{\notin i}\alpha + \mathbf{U}_i\beta \ ,$$

it's possible to solve for $\alpha$ and for $\beta$ and then retrieve the solution of the underdetermined linear system. Projecting the linear system on $\mathbf{V}_{\notin i}$, and prescibing the orthogonality condition $\mathbf{V}_i^*\mathbf{B}\widetilde{\mathbf{u}_{i/p}} = \mathbf{0}$ on the solution, two decoupled sets of equations for $\alpha$ and $\beta$ appears,

$$\begin{cases} \mathbf{V}_{\notin i}^* \left( \mathbf{A} - s_i\mathbf{B} \right) \mathbf{U}_{\notin i}\alpha = \mathbf{V}_{\notin i}^*\mathbf{b}_i \\ \mathbf{V}_i^*\mathbf{B}\mathbf{U}_i\beta = \mathbf{0} \ , \end{cases}$$

as $(\mathbf{A} - s_i\mathbf{B}) \mathbf{U}_i = \mathbf{0}$ and $\mathbf{V}_i^*\mathbf{B}\mathbf{U}_{\notin i} = \mathbf{0}$. Exploiting diagonalization properties, the system can be written as two uncoupled diagonal systems

$$\begin{cases} \mathrm{diag} \left\{ a^{(j)\notin i} - s_i b^{(j)\notin i} \right\} \alpha = \mathbf{V}_{\notin i}^*\mathbf{b}_i \\ \mathrm{diag} \left\{ b^{(i)} \right\} \beta = \mathbf{0} \ , \end{cases}$$

---

**5.2. Sensitivity of spectral decomposition**

or, exploiting normalization condition $\mathbf{v}_i^* \mathbf{B} \mathbf{u}_i = 1$, and thus $b^{(j)} = 1$ and $a^{(j)} = s_j$,

$$\begin{cases} \text{diag}\left\{ s_{j \neq i} - s_i \right\} \alpha = \mathbf{V}_{\neq i}^* \mathbf{b}_i \\ \mathbf{I}\,\beta = \mathbf{0}\,, \end{cases}$$

The solution reads

$$\begin{cases} \alpha = \text{diag}\left\{ \dfrac{1}{s_{j \neq i} - s_i} \right\} \mathbf{V}_{\neq i}^* \mathbf{b}_i \\ \beta = \mathbf{0}\,, \end{cases}$$

and thus

$$\tilde{\mathbf{u}}_{i/p} = \mathbf{U}_{\neq i} \text{diag}\left\{ \frac{1}{s_{j \neq i} - s_i} \right\} \mathbf{V}_{\neq i}^* \mathbf{b}_i =$$

$$= \mathbf{U}_{\neq i} \text{diag}\left\{ \frac{1}{s_{j \neq i} - s_i} \right\} \mathbf{V}_{\neq i}^* (-(\mathbf{A}_{/p} - s_{i/p}\mathbf{B} - s_i \mathbf{B}_{/p})\mathbf{u}_i) =$$

$$= \mathbf{U}_{\neq i} \text{diag}\left\{ \frac{1}{s_{j \neq i} - s_i} \right\} \mathbf{V}_{\neq i}^* (-(\mathbf{A}_{/p} - s_i \mathbf{B}_{/p})\mathbf{u}_i) =$$

$$= \sum_{j \neq i} \frac{\mathbf{v}_j^* (-\mathbf{A}_{/p} + s_i \mathbf{B}_{/p})\mathbf{u}_i}{s_j - s_i} \mathbf{u}_j\,.$$

### Some algebra with components

$$u_{a/p}^i = \sum_{b \neq i} u_a^{(b)} \delta_{bc} D^{(c)} v_d^{(c)\,*} b_d^{(i)} =$$

$$= \sum_{b \neq i} u_a^{(b)} \delta_{bc} D^{(c)} v_d^{(c)\,*} M_{de} u_e^{(i)} =$$

$$= \sum_{b \neq i} u_a^{(b)} D^{(b)} v_d^{(b)\,*} M_{de} u_e^{(i)} =$$

$$= \sum_{b \neq i} C^{b,i} u_a^{(b)}\,,$$

with $C^{b,i} = D^{(b)} v_d^{(b)\,*} M_{de/p} u_e^{(i)} = \frac{\mathbf{v}_b^* \mathbf{M} \mathbf{u}_i}{s_b - s_i}$, and $\mathbf{M} = -\mathbf{A}_{/p} + s_i \mathbf{B}_{/p}$.

## 5.2.2 Generalized eigenvalue problem (second order)

### Right and left eigenvalue problem

# Part II

# Multivariable Calculus

# INTRODUCTION TO MULTI-VARIABLE CALCULUS

## 6.1 Function

## 6.2 Limit

## 6.3 Derivatives

## 6.4 Integrals

## 6.5 Theorems

### 6.5.1 Green's lemma

$$\int_S \frac{\partial F}{\partial y} dx dy = - \oint_{\partial S} F dx$$
$$\int_S \frac{\partial G}{\partial x} dx dy = \oint_{\partial S} G dy$$

**Proof for simple domains.**

In a simple domain in $x$, so that the closed contour $\partial S$ is delimited by the curves $y = Y_1(x)$, $y = Y_2(x) > Y_1(x)$, for $x \in [x_1, x_2]$,

$$\int_S \frac{\partial F}{\partial y} dx dy = \int_{x=x_1}^{x_2} \int_{y=Y_1(x)}^{Y_2(x)} \frac{\partial F}{\partial y} dy \, dx =$$
$$= \int_{x=x_1}^{x_2} \left[ F(x, Y_2(x)) - F(x, Y_1(x)) \right] dx =$$
$$= - \int_{x=x_1}^{x_2} F(x, Y_1(x)) - \int_{x=x_2}^{x_1} F(x, Y_2(x)) dx =$$
$$= - \oint_{\partial S} F(x, y) dx$$

## Proof for generic domain

**todo**

In a simple domain in $y$, so that the closed contour $\partial S$ is delimited by the curves $x = X_1(y)$, $x = X_2(y) > X_1(y)$ for $y \in [y_1, y_2]$,

$$\int_S \frac{\partial G}{\partial x} dx dy = \int_{y=y_1}^{y_2} \int_{x=X_1(y)}^{X_2(y)} \frac{\partial G}{\partial x} dx\, dy =$$

$$= \int_{y=y_1}^{y_2} [G(X_2(y), y) - G(X_1(y), y)]\, dy =$$

$$= \int_{y=y_1}^{y_2} G(X_1(y), y) dy + \int_{y=y_2}^{y_1} G(X_2(y), y) dy =$$

$$= \oint_{\partial S} G(x, y) dy$$

# Part III

# Differential Geometry

# INTRODUCTION TO DIFFERENTIAL GEOMETRY

## 7.1 Differential geometry in $E^3$

### 7.1.1 Curves

Parametric representation of curve in 3-dimensional (Euclidean) space $E^3$

$$\vec{r}(q)$$

**Differential, $d\vec{r}$.**

$$d\vec{r}(q) = \vec{r}'(q)\,dq\;.$$

**Arc-length parameter, $s$.** So that $ds = |d\vec{r}(s)|$ and thus

$$|d\vec{r}(s)| = |\vec{r}'(s)|\,|ds| \qquad \rightarrow \qquad |\vec{r}'(s)| = 1 \qquad \rightarrow \qquad \vec{r}'(s) = \hat{t}(s)\;.$$

**Frenet basis.** Using arc-length parameter, Frenet basis is naturally defined as the set $\{\hat{t}, \hat{n}, \hat{b}\}$:

- tangent unit vector, $\hat{t}(s) = \vec{r}'(s)$,
- normal unit vector, $\hat{r}''(s) = \hat{t}'(s) =: \kappa(s)\,\hat{n}(s)$, with $\kappa(s)$ local curvature
- binormal unit vector, $\hat{b}(s) = \hat{t}(s) \times \hat{n}(s)$

Using a general parameter, $t$, with some abuse of notation $\vec{r}(t) = \vec{r}(s(t))$ and indicating $\dot{(\,)} = \frac{d}{dt}$,

- $\dot{\vec{r}} = \frac{ds}{dt}\frac{d\vec{r}}{ds} = \dot{s}\hat{t}$
- $\ddot{\vec{r}} = \frac{d}{dt}\dot{\vec{r}} = \frac{d}{dt}\left(\dot{s}\hat{t}\right) = \ddot{s}\hat{t} + \frac{ds}{dt}\frac{d}{ds}\hat{t} = \ddot{s}\hat{t} + \dot{s}^2\kappa\,\hat{n}$

**Osculator circle.** Circle with $R(s) = \frac{1}{\kappa(s)}$, in plane orthogonal to $\hat{b}(s)$, passing through $\vec{r}(s)$, and thus center in $\vec{r}_C(s) = \vec{r}(s) + \hat{n}R(s)$. Its parametric representation using its arc-length parameter $p$, with $\vec{r}(p=0) = \vec{r}(s)$ reads

$$\vec{r}(p) = \vec{r}_C(s) + R(s)\left[-\cos\left(\frac{p}{R(s)}\right)\hat{n}(s) + \sin\left(\frac{p}{R(s)}\right)\hat{t}(s)\right]\;.$$

Its first and second order derivatives w.r.t. the arc-length $p$ evaluated in $p = 0$, i.e. $\vec{r} = \vec{r}(s)$ read:

- first derivative in $p = 0$,

$$\hat{t}(p)\big|_{p=0} = \vec{r}'(p)\big|_{p=0} = \left[\sin\left(\frac{p}{R(s)}\right)\hat{n}(s) + \cos\left(\frac{p}{R(s)}\right)\hat{t}(s)\right]\Bigg|_{p=0} = \hat{t}(s)\;,$$

i.e. the osculator circle has the same tangent as the curve in the point.

- second derivative in $p = 0$,

$$\kappa(p)\hat{n}(p)\big|_{p=0} = \vec{r}''(p)\big|_{p=0} = \frac{1}{R(s)}\left[\cos\left(\frac{p}{R(s)}\right)\hat{n}(s) - \sin\left(\frac{p}{R(s)}\right)\hat{t}(s)\right]\bigg|_{p=0} = \frac{1}{R(s)}\hat{n}(s) = \kappa(s)\hat{n}(s)\,,$$

i.e. the osculator circle has the same normal vector and curvature as the curve in the point.

## 7.1.2 Surfaces

$$\vec{r}(q^1, q^2)$$

$$d\vec{r} = \frac{\partial\vec{r}}{\partial q^1}\,dq^1 + \frac{\partial\vec{r}}{\partial q^2}\,dq^2 = \vec{b}_1\,dq^1 + \vec{b}_2\,dq^2$$

A third vector $\vec{b}_3 := \hat{n}$ can be defined so that $|\hat{n}| = 1$ and $\hat{n}\cdot\vec{b}_i = 0$, $i = 1:2$. For $i = 1:2$, $k = 1:2$

$$\frac{\partial\vec{b}_i}{\partial q^j} = \Gamma_{ij}^k\vec{b}_k = \Gamma_{ij}^1\vec{b}_1 + \Gamma_{ij}^2\vec{b}_2 + \Gamma_{ij}^3\vec{b}_3$$

so that

$$\Gamma_{ij}^k = \vec{b}^k\cdot\frac{\partial\vec{b}_i}{\partial q^j}$$

**Normal vector.**

$$\vec{n}(q^1, q^2) = \frac{\partial\vec{r}}{\partial q^1}(q^1, q^2) \times \frac{\partial\vec{r}}{\partial q^2}(q^1, q^2) = \vec{b}_1(q^1, q^2) \times \vec{b}_2(q^1, q^2)$$

**Tangent plane.**

$$(\vec{r} - \vec{r}(q^1, q^2))\cdot\vec{n}(q^1, q^2) = 0$$

**Length of elementary segment.**

$$|d\vec{r}|^2 = d\vec{r}\cdot d\vec{r} =$$
$$= \left(\vec{b}_1\,dq^1 + \vec{b}_2\,dq^2\right)\cdot\left(\vec{b}_1\,dq^1 + \vec{b}_2\,dq^2\right) =$$
$$= g_{11}\,dq^1\,dq^1 + g_{12}\,dq^1\,dq^2 + g_{21}\,dq^2\,dq^1 + g_{22}\,dq^2\,dq^2 = g_{ij}\,dq^i\,dq^j$$

**Second order approximation.**

$$\vec{r}(q^1 + dq^1, q^2 + dq^2) = \vec{r}(q_1, q_2) + \frac{\partial\vec{r}}{\partial q^i}\,dq^i + \frac{\partial^2\vec{r}}{\partial q^i\partial q^j}\,dq^i\,dq^j =$$
$$= \vec{r}(q_1, q_2) + \vec{b}_i\,dq^i + \vec{b}_k\Gamma_{ij}^k\,dq^i\,dq^j + \hat{n}\,\Gamma_{ij}^3\,dq^i\,dq^j$$

so that

$$\left[\vec{r}(q^1 + dq^1, q^2 + dq^2) - \vec{r}(q^1, q^2)\right]\cdot\hat{n} = \Gamma_{ij}^3\,dq^i\,dq^j =$$
$$= \hat{n}\cdot\frac{\partial^2\vec{r}}{\partial q^i\partial q^j}\,dq^i\,dq^j =$$
$$= \hat{n}\cdot\frac{\partial^2\vec{r}}{\partial q^i\partial q^j}\,\vec{b}^i\cdot\vec{b}_k dq^k\,\vec{b}^j\cdot\vec{b}_l dq^l =$$
$$= \underbrace{dq^k\vec{b}_k}_{d\vec{r}}\cdot\left[\hat{n}\cdot\frac{\partial^2\vec{r}}{\partial q^i\partial q^j}\vec{b}^i\otimes\vec{b}^j\right]\cdot\underbrace{dq^l\vec{b}_l}_{d\vec{r}}$$

**Curvature tensor.**

# Part IV

# Vector and Tensor Algebra and Calculus

# TENSOR ALGEBRA

This section introduces tensor algebra

> **Warning:** This introduction is meant for **spaces with inner product** that allows to introduce quite naturally the useful concept of the *reciprocal basis* of a given basis of $\mathcal{V}$, belonging to the same space $\mathcal{V}$, somehow avoiding the complications coming from the introduction of dual space $\mathcal{V}^*$ and basis, and everything is required for that.

## 8.1 Vector space $\mathcal{V}$

A vector space is an algebraic structure with:

- a **set** $\mathcal{V}$, whose elements are called **vectors**, here indicated with bold symbols $\mathbf{v} \in \mathcal{V}$

- a **field** $K$ (usually $\mathbb{R}$ or $\mathbb{C}$), whose elements are calles **scalars**,

- 2 **operations**, closed w.r.t. the set $\mathcal{V}$:

    - **vector sum**, $\mathbf{u} + \mathbf{v} \in \mathcal{V}$ for $\forall \mathbf{u},\ \mathbf{v} \in \mathcal{V}$

    - **multiplication of a scalar and a vector**, $a\mathbf{v} \in \mathcal{V}$ for $\forall a \in K,\ \mathbf{u} \in \mathcal{V}$

    with properties discussed below **todo** …

### 8.1.1 Operations (I)

**Sum**

…

**Multiplication by a scalar**

…

---

**Definition 8.1.1 (Linear combination)**

A linear combination of $n$ vectors $\{\mathbf{v}_i\}_{i=1:n}$, $\mathbf{v}_i \in \mathcal{V}$, is the weighted sum

$$a^1\mathbf{v}_1 + a^2\mathbf{v}_2 + \cdots + a^n\mathbf{v}_n = a^i\mathbf{v}_i \in \mathcal{V}\,,$$

---

having used Einstein's summation convention over repeated indices. Here the position of the indices has no particular meaning, but it'll have soon in the following sections.

## Inner product

**The existence of an inner product is not a requirement of a vector space**

…

An inner product in a vectors space $\mathcal{V}$ over field $K$ is an operation $\cdot : \mathcal{V} \times \mathcal{V} \to K$,

$$\mathbf{u} \cdot \mathbf{v} \, ,$$

with the following properties: **todo**

## 8.1.2 Basis of a vector space $\mathcal{V}$

**Definition 8.1.2 (Basis)**

A basis is a minimal set of vectors of $\mathcal{V}$ that can represent all the elements of $\mathcal{V}$.

**Definition 8.1.3 (Dimension of a vector space)**

The dimension of a vector space $\mathcal{V}$ is the number of the elements of a basis of the space.

**Definition 8.1.4 (Reciprocal basis)**

In a inner product space, the reciprocal basis of a given basis $\{\mathbf{b}_a\}_{a=1:d}$ is the set of vectors $\{\mathbf{b}^b\}_{b=1:d}$, s.t.

$$\mathbf{b}^b \cdot \mathbf{b}_a = \delta^b_a \, .$$

**Definition 8.1.5 ("Metric tensor")**

$$g_{ab} := \mathbf{b}_a \cdot \mathbf{b}_b$$
$$g^{ab} := \mathbf{b}^a \cdot \mathbf{b}^b$$

The following holds

$$\begin{aligned} \mathbf{b}_a &= g_{ab}\mathbf{b}^b \quad (1) \\ \mathbf{b}^a &= g^{ab}\mathbf{b}_b \quad (2) \end{aligned} \tag{8.1}$$

### "Proof"

Taking the dot product with $\mathbf{b}_c$ of the relation (8.1)(1),

$$\mathbf{b}_c \cdot \mathbf{b}_a = g_{ab} \underbrace{\mathbf{b}_c \cdot \mathbf{b}^b}_{=\delta_c^b} = g_{ac} .$$

## 8.1.3 Change of basis

Let $T_a^b$ the elements of the matrix of change of basis, representing the vectors of the basis $\{\tilde{b}_a\}$ as linear combinations of the vectors of the basis $\{\mathbf{b}_b\}$,

$$\tilde{\mathbf{b}}_b = T_b^a \mathbf{b}_a .$$

**Inverse transformation.** Let $\widetilde{T}$ be the elements of the inverse transformation,

$$\mathbf{b}_c = \widetilde{T}_c^b \tilde{\mathbf{b}}_b = \widetilde{T}_c^b T_b^a \mathbf{b}_a ,$$

and thus

$$\widetilde{T}_c^b T_b^a = \delta_c^a .$$

**Transformation of the reciprocal basis. todo**

### Transformation of components

A vector $\mathbf{v}$ can be represented in different basis, as different linear combinations of the elements of those bases,

$$\mathbf{v} = v^a \mathbf{b}_a = \tilde{v}^b \tilde{\mathbf{b}}_b .$$

Given the rules of change of basis, the rule of transformation of components immediately follwos

$$\tilde{\mathbf{b}}_b = T_b^a \mathbf{b}_a \quad \tilde{v}^b = \widetilde{T}_a^b v^a$$
$$\mathbf{b}_b = \widetilde{T}_b^a \tilde{\mathbf{b}}_a \quad v^b = T_a^b \tilde{v}^a$$

### Proof

$$\mathbf{v} = v^a \mathbf{b}_a = \underbrace{v^a \widetilde{T}_a^b}_{\tilde{v}^b} \tilde{\mathbf{b}}_b = \tilde{v}^b \tilde{\mathbf{b}}_b$$

or

$$\mathbf{v} = \tilde{v}^b \tilde{\mathbf{b}}_b = \underbrace{\tilde{v}^b \widetilde{T}_c^b}_{v^c} \mathbf{b}_c = v^c \mathbf{b}_c$$

It's clear that the vectors of the bases and the components follow inverse transformations to preserve the **invariance of vector w.r.t. a change of basis**: the vector $\mathbf{v}$ doesn't change if we change our description of it, by changing a basis.

---

**Matrix of change of basis as a tensor (todo maybe later? Tensor not introduced yet here)**

The rule of transformation between different basis can be interpreted using dot product between tensors and vectors

$$\tilde{\mathbf{b}}_b = T_b^a \mathbf{b}_a = T_c^a \mathbf{b}_a \delta_b^c = T_c^a \mathbf{b}_a \otimes \mathbf{b}^c \cdot \mathbf{b}_b = (T_c^a \mathbf{b}_a \otimes \mathbf{b}^c) \cdot \mathbf{b}_b = \mathbb{T} \cdot \mathbf{b}_b$$

**Inverse and transpose**

Interpreting the indices of the transformation matrix as the indices of rows and columns of a 2x2 matrix, transformation of components involves the transpose of the inverse matrix, as the indices $a$, $b$ are swapped.

**Example 8.1.1**

**Example 8.1.2 (Rotation)**

As the inverse of a rotation is its transpose, $\widetilde{T}_a^b = T_b^a$. So the rule of transformation of components follows the same rule of change of basis. As an example, let the transformation between two Cartesian bases be the rotation

$$\begin{aligned} \tilde{\mathbf{x}} &= \cos\theta\,\mathbf{x} + \sin\theta\,\mathbf{y} \\ \tilde{\mathbf{y}} &= -\sin\theta\,\mathbf{x} + \cos\theta\,\mathbf{y} \end{aligned} \quad, \quad \begin{aligned} \mathbf{x} &= \cos\theta\,\tilde{\mathbf{x}} - \sin\theta\,\tilde{\mathbf{y}} \\ \mathbf{y} &= \sin\theta\,\tilde{\mathbf{x}} + \cos\theta\,\tilde{\mathbf{y}} \end{aligned}$$

Let a vector

$$\begin{aligned} \mathbf{v} &= v_x\mathbf{x} + v_y\mathbf{y} = \\ &= v_x\left(\cos\theta\,\tilde{\mathbf{x}} - \sin\theta\,\tilde{\mathbf{y}}\right) + v_y\left(\sin\theta\,\tilde{\mathbf{x}} + \cos\theta\,\tilde{\mathbf{y}}\right) = \\ &= \tilde{v}_x\tilde{\mathbf{x}} + \tilde{v}_y\tilde{\mathbf{y}} = \\ &= \tilde{v}_x\left(\cos\theta\,\mathbf{x} + \sin\theta\,\mathbf{y}\right) + \tilde{v}_y\left(-\sin\theta\,\mathbf{x} + \cos\theta\,\mathbf{y}\right) \end{aligned}$$

and thus

$$\begin{aligned} \tilde{v}_x &= \cos\theta\,v_x + \sin\theta\,v_y \\ \tilde{v}_y &= -\sin\theta\,v_x + \sin\theta\,v_y \end{aligned} \quad, \quad \begin{aligned} v_x &= \cos\theta\,\tilde{v}_x - \sin\theta\,\tilde{v}_y \\ v_y &= \sin\theta\,\tilde{v}_x + \sin\theta\,\tilde{v}_y \end{aligned}$$

## 8.1.4 Operations

**Tensor product of vectors**

$$\mathbf{v}_{(1)} \otimes \mathbf{v}_{(2)} \otimes \cdots \otimes \mathbf{v}_{(r)}$$

or writing each vector as a linear combination of the elements of a basis $\mathbf{b}_a$,

$$\mathbf{v}_{(i)} = v_{(i)}^{a_i}\mathbf{b}_{a_i}\,, \cdots$$

$$\begin{aligned} \mathbf{v}_{(1)} \otimes \mathbf{v}_{(2)} \otimes \cdots \otimes \mathbf{v}_{(3)} &= \left(v_{(1)}^{a_1}\mathbf{b}_{a_1}\right) \otimes \left(v_{(2)}^{a_2}\mathbf{b}_{a_2}\right) \otimes \cdots \otimes \left(v_{(r)}^{a_r}\mathbf{b}_{a_r}\right) = \\ &= v_{(1)}^{a_1}\,v_{(2)}^{a_2}\,\cdots\,v_{(r)}^{a_r}\mathbf{b}_{a_1} \otimes \mathbf{b}_{a_2} \otimes \cdots \otimes \mathbf{b}_{a_r} \end{aligned}$$

The result of the tensor product of $r$ vectors is a rank-$r$ tensor, $\mathbb{T}$, as it will be clear below, with components

$$T^{a_1 a_2 \cdots a_r} = v_{(1)}^{a_1}\,v_{(2)}^{a_2}\,\cdots\,v_{(r)}^{a_r}\,.$$

## 8.2 Space of tensors

---

**Definition 8.2.1 (Tensor as a multilinear map)**

A rank-$r$ tensor $\mathbf{A}$ can be defined as a multilinear map, acting on $r$ vectors[1] $\mathbf{v}^{(i)} \in \mathcal{V}$, $i = 1 : r$

$$\mathbf{A}\left(\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(r)}\right) : \underbrace{\mathcal{V} \times \cdots \times \mathcal{V}}_{r \text{ times}} \to K.$$

---

**todo**

- Tensor definition by components, and action on the vectors

- Reference to a similar definition of a vector: a vector $\mathbf{v} \in \mathcal{V}$, vector space with inner product, is uniquely determined by the value of its scalar product $\mathbf{u} \cdot \mathbf{v}$ with every vector $\mathbf{u} \in \mathcal{V}$ or, analogously, with every vector in a set of $n$ linearly independent vectors $\left\{\mathbf{u}^{(i)}\right\}$, $i = 1 : n$, with $n = \dim(\mathcal{V})$.

- Definition of components w.r.t. a basis

## 8.2.1 Operations (I)

### Sum

Sum of tensors with the same rank, $\mathbf{A}, \mathbf{B} \in \mathcal{V}$,

$$\mathbf{A} + \mathbf{B} \in \mathcal{V}^r$$

**todo** Components

### Multiplication by a scalar

Multiplication by a scalar $a \in K$, of a $r$-rank tensor $\mathbf{A}$,

$$a\mathbf{A} \in \mathcal{V}^r$$

**todo** Components

---

**Vector space of tensors**

The set of tensors of a given rank with the operations of sum and multiplication by a scalar defined above forms a vector space.

---

[1] This introduction has no ambition iof being the most general — and precise? — introduction to tensors: as already stated somewhere else, existence of inner product is assumed even it's not necessary for the most general treatment of tensors; moreover, no distinction is made between dual space $\mathcal{V}^*$ of covectors — space of linear forms on $V$ — and $V$ itself;...

## 8.2.2 Basis

### Change of basis and rule of transformation of components - classical definition of a tensor

$$\tilde{\mathbf{b}}_{i_a} = T^{i_a}_{i_b}\mathbf{b}_{i_b}$$

$$\mathbf{b}_{i_a} = \widetilde{T}^{i_b}_{i_a}\tilde{\mathbf{b}}_{i_b}$$

$$\mathbf{A} = A^{i_1\ldots i_p}\mathbf{b}_{i_1}\ldots\mathbf{b}_{i_p} =$$

$$= A^{i_1\ldots i_p}\left(T^{j_1}_{i_1}\tilde{\mathbf{b}}_{j_1}\right)\ldots\left(T^{j_p}_{i_p}\tilde{\mathbf{b}}_{j_p}\right) =$$

$$= A^{i_1\ldots i_p}T^{j_1}_{i_1}\ldots T^{j_p}_{i_p}\tilde{\mathbf{b}}_{j_1}\ldots\tilde{\mathbf{b}}_{j_p} =$$

$$= \widetilde{A}^{j_1\ldots j_p}\tilde{\mathbf{b}}_{j_1}\ldots\tilde{\mathbf{b}}_{j_p}\,,$$

with

$$\widetilde{A}^{j_1\ldots j_p} = A^{i_1\ldots i_p}T^{j_1}_{i_1}\ldots T^{j_p}_{i_p}\,. \tag{8.2}$$

---

**Definition 8.2.2 (Classical definition of a tensor)**

Relation (8.2) is the "historical" definition of a tensor, through the law of transformation of its components following a change of basis.

---

## 8.2.3 Operations (II)

### Tensor product

The tensor product of a $p$-rank tensor $\mathbf{A}$ and a $q$-rank tensor $\mathbf{B}$ is the $(p+q)$-rank tensor $\mathbf{A}\otimes\mathbf{B} = \mathbf{AB}$, that can be defined using component representation in a given basis,

$$\mathbf{A}\otimes\mathbf{B} = \left(A^{a_1\ldots a_p}\mathbf{b}_{a_1}\ldots\mathbf{b}_{a_p}\right)\otimes\left(B^{b_1\ldots b_q}\mathbf{b}_{b_1}\ldots\mathbf{b}_{b_q}\right) =$$

$$= A^{a_1\ldots a_{p-1}a_p}B^{b_1\,b_2\ldots b_q}\mathbf{b}_{a_1}\ldots\mathbf{b}_{a_p}\mathbf{b}_{b_1}\ldots\mathbf{b}_{b_q}$$

### Dot product

The dot product of a $p$-rank tensor $\mathbf{A}$ and a $q$-rank tensor $\mathbf{B}$ is the $(p+q-2)$-rank tensor $\mathbf{A}\cdot\mathbf{B}$, that can be defined using component representation in a given basis,

$$\mathbf{A}\cdot\mathbf{B} = \left(A^{a_1\ldots a_p}\mathbf{b}_{a_1}\ldots\mathbf{b}_{a_p}\right)\cdot\left(B_{b_1\ldots b_q}\mathbf{b}^{b_1}\ldots\mathbf{b}^{b_q}\right) =$$

$$= A^{a_1\ldots a_p}B_{b_1\ldots b_q}\mathbf{b}_{a_1}\ldots\underbrace{\mathbf{b}_{a_p}\cdot\mathbf{b}^{b_1}}_{=\delta^{b_1}_{a_p}}\ldots\mathbf{b}^{b_q} =$$

$$= A^{a_1\ldots a_{p-1}k}B_{k\,b_2\ldots b_q}\mathbf{b}_{a_1}\ldots\mathbf{b}_{a_{p-1}}\mathbf{b}^{b_2}\ldots\mathbf{b}^{b_q}$$

**Contraction**

Contraction of a pair of index of a $p$-rank tensor $\mathbf{A}$ returns a $p-2$-rank tensor defined as

$$C_a^b\left(\mathbf{A}\right) = ...$$

**Exterior product**

**todo** *see exterior algebra*

## 8.2.4 Invariants

# 8.3 Exterior algebra

$$\wedge$$

## 8.3.1 Exterior product

Generalization of the vector product

# TENSOR CALCULUS IN EUCLIDEAN SPACES

This section deals with tensor calculus in Euclidean space or on manifolds embedded in Euclidean spaces, focusing on $d$-dimensional spaces with $d \leq 3$, with *inner product*.

This section may rely on results of *differential geometry*.

## 9.1 Coordinates

A set of parameters $\{q^a\}_{a=1:d}$ to represent vector (or point) in space,

$$\vec{r}(q^a)$$

if $\vec{r} \in E^d$, $a = 1 : d$.

In $E^3$,

- **Coordinate lines**, 2-parameter family of lines, keeping 2 coordinates constant. As an example, coordinate lines with constant $q^2$, $q^3$

$$\vec{r}_1(q^1) = \vec{r}(q^1, \bar{q}^2, \bar{q}^3) \ .$$

- **Coordinate surfaces,** 1-parameter family of surfaces, keeping 1 coordinate constant. As an example, coordinate surfaces with constant $q^1$,

$$\vec{r}_{23}(q^2, q^3) = \vec{r}(\bar{q}^1, q^2, q^3) \ .$$

**Definition 9.1.1 (Regular parametrization)**

If $\frac{\partial \vec{r}}{\partial q^a} \neq 0$.

### 9.1.1 Natural basis

**Definition 9.1.2 (Natural basis)**

Vectors of natural basis

$$\vec{b}_a := \frac{\partial \vec{r}}{\partial q^a}$$

**Definition 9.1.3 (Reciprocal basis (todo move to Tensor Algebra))**

Given a basis $\{\vec{b}_a\}_a$, its reciprocal basis the set of vector $\{\vec{b}^b\}_b$ defined as

$$\vec{b}^b \cdot \vec{b}_a = \delta_a^b \ ,$$

being $\delta_a^b$ Kronecker delta.

**Definition 9.1.4 (Christoffel symbols)**

Christoffel symbols (of the $2^{nd}$ kind) are defined as the components of the partial derivatives of the vectors of a natural basis w.r.t. the coordinates referred to the natural basis itself,

$$\frac{\partial \vec{b}_a}{\partial q^b} = \Gamma_{ab}^c \, \vec{b}_c \tag{9.1}$$

### Properties of Christoffel symbols

Exploiting the definition of reciprocal basis, Christoffel symbols can be written as

$$\Gamma_{ab}^c = \vec{b}^c \cdot \frac{\partial \vec{b}_a}{\partial q^b} \ .$$

**Symmetry.** Symmetry os the lower indices

$$\Gamma_{ab}^c = \Gamma_{ba}^c \ ,$$

readily follows Schwartz theorem about partial derivatives

$$\frac{\partial \vec{b}_a}{\partial q^b} = \frac{\partial}{\partial q^c} \frac{\partial \vec{r}}{\partial q^a} = \frac{\partial}{\partial q^a} \frac{\partial \vec{r}}{\partial q^b} = \frac{\partial \vec{b}_b}{\partial q^a}$$

## 9.2 Fields

Function of the points in space $F : E^d \to V^r$, being $V^r$ a space of tensors of order $r$.

## 9.3 Differential operators

### 9.3.1 Directional derivative

$$F(\vec{r}) = F\left(\vec{r}\left(q^a\right)\right) = f(q^a)$$

$$f(q^a + \beta \Delta q^a) = F(\vec{r}(q^a + \beta \Delta q^a))$$

$$\vec{r}(q^a) + \alpha \vec{v} = \vec{r}(q^a + \beta \Delta q^a) \sim \vec{r}(q^a) + \frac{\partial \vec{r}}{\partial q^b} \beta \Delta q^b$$

$$\alpha \vec{v} \sim \beta \frac{\partial \vec{r}}{\partial q^b}(q^a) \, \Delta q^b = \beta \vec{b}_b(q^a) \Delta q^b \qquad \to \qquad \Delta q^b = \frac{\alpha}{\beta} \vec{b}^b(q^a) \cdot \vec{v}$$

The directional derivative for an arbitrary vector $\vec{v} \in V$

$$\left. \frac{d}{d\alpha} F(\vec{r} + \alpha\vec{v}) \right|_{\alpha=0}$$

is evaluated as the limit for $\alpha \to 0$ of the incremental ratio

$$\frac{F(\vec{r} + \alpha\vec{v}) - F(\vec{r})}{\alpha} \sim \frac{f(q^a + \beta\Delta q^a) - f(q^a)}{\alpha} =$$

$$\sim \frac{1}{\alpha} \frac{\partial f}{\partial q^b}(q^a)\beta\Delta q^b =$$

$$\sim \vec{v} \cdot \vec{b}^b(q^a)\frac{\partial f}{\partial q^b}(q^a) =$$

$$= \vec{v} \cdot \nabla F(\vec{r})$$

### 9.3.2 Gradient

The gradient is the differential operator is the first-order differential operator appearing in the definition of the directional derivative, $\nabla F(\vec{r})$. It takes a tensor field $F(\vec{r})$ of order $r$ and gives a tensor field $\nabla F(\vec{r})$ of order $r + 1$. Given a set of coordinates $\{q^a\}_{a=1:d}$, the gradient can be written using the reciprocal basis of the natural basis as

$$\nabla F(\vec{r}) = \vec{b}^b(\vec{r})\frac{\partial F}{\partial q^b}(\vec{r}) \tag{9.2}$$

**Examples.** …

---

**Example 9.3.1 (Gradient of a scalar field - with general coordinates $q^a$)**

Applying the definition (9.2) of gradient operator, it readily follows

$$\nabla F = \vec{b}^a \frac{\partial F}{\partial q^a}$$

---

**Example 9.3.2 (Gradient of a vector field - with general coordinates $q^a$)**

Applying the definition (9.2) of gradient operator, rule for the derivative of a product and the definition (9.1) of Christoffel symbols to write derivatives of base vectors,

$$\nabla F = \vec{b}^a \frac{\partial}{\partial q^a}\left(F^b\vec{b}_b\right) =$$

$$= \vec{b}^a \left[\frac{\partial F^b}{\partial q^a}\vec{b}_b + F^b \frac{\partial \vec{b}_b}{\partial q^a}\right] =$$

$$= \vec{b}^a \left[\frac{\partial F^b}{\partial q^a}\vec{b}_b + F^b \Gamma^c_{ab}\vec{b}_c\right] =$$

$$= \vec{b}^a \otimes \vec{b}_b \left[\frac{\partial F^b}{\partial q^a} + \Gamma^b_{ac}F^c\right] \ .$$

---

**Example 9.3.3 (Gradient of a $2^{nd}$-order tensor field - with general coordinates $q^a$)**

Applying the definition (9.2) of gradient operator, rule for the derivative of a product and the definition (9.1) of Christoffel symbols to write derivatives of base vectors,

$$\nabla F = \vec{b}^a \frac{\partial}{\partial q^a} \left( F^{bc} \vec{b}_b \otimes \vec{b}_c \right) =$$

$$= \vec{b}^a \left[ \frac{\partial F^{bc}}{\partial q^a} \vec{b}_b \vec{b}_c + F^{bc} \frac{\partial \vec{b}_b}{\partial q^a} \vec{b}_c + F^{bc} \vec{b}_b \frac{\partial \vec{b}_c}{\partial q^a} \right] =$$

$$= \vec{b}^a \left[ \frac{\partial F^{bc}}{\partial q^a} \vec{b}_b \vec{b}_c + F^{bc} \Gamma^d_{ab} \vec{b}_d \vec{b}_c + F^{bc} \Gamma^{ac}_d \vec{b}_b \vec{b}_d \right] =$$

$$= \vec{b}^a \otimes \vec{b}_b \otimes \vec{b}_c \left[ \frac{\partial F^{bc}}{\partial q^a} + \Gamma^b_{ad} F^{dc} + \Gamma^c_{ad} F^{bd} \right] \ .$$

### 9.3.3 Divergence

Divergence opearator is a first-order differential operator that can be defined as the contraction of the first two indices of the gradient,

$$\nabla \cdot F = C^2_1 \left( \nabla F \right) \ .$$

It takes a tensor field $F(\vec{r})$ of order $r \geq 1$ and gives a tensor field $\nabla \cdot F(\vec{r})$ of order $r - 1 \geq 0$.

**Example 9.3.4 (Divergence of a vector field - with general coordiantes $q^a$)**

Applying contraction to the gradient of a vector field, it readily follows,

$$\nabla \cdot \left( F^b \vec{b}_b \right) = C^2_1 \left( \nabla F \right) =$$

$$= C^2_1 \left( \vec{b}^a \otimes \vec{b}_b \left[ \frac{\partial F^b}{\partial q^a} + \Gamma^b_{ac} F^c \right] \right) =$$

$$= \frac{\partial F^a}{\partial q^a} + \Gamma^a_{ac} F^c$$

**Example 9.3.5 (Divergence of a $2^{nd}$-order tensor field - with general coordiantes $q^a$)**

Applying contraction to the gradient of a vector field, it readily follows,

$$\nabla \cdot \left( F^{bc} \vec{b}_b \otimes \vec{b}_c \right) = C^2_1 \left( \nabla F \right) =$$

$$= C^2_1 \left( \vec{b}^a \otimes \vec{b}_b \otimes \vec{b}_c \left[ \frac{\partial F^{bc}}{\partial q^a} + \Gamma^b_{ad} F^{dc} + \Gamma^c_{ad} F^{bd} \right] \right) =$$

$$= \vec{b}_c \left[ \frac{\partial F^{ac}}{\partial q^a} + \Gamma^a_{ad} F^{dc} + \Gamma^c_{ad} F^{ad} \right]$$

## 9.3.4 Laplacian

Laplacian operator is second-order differential operator that can be defined as the divergence of the gradient,

$$\Delta F = \nabla^2 F = \nabla \cdot \nabla F \ .$$

---

**Example 9.3.6 (Laplacian of a scalar field - with general coordinates $q^a$)**

$$\nabla \cdot \nabla F = C_1^2 \left[ \nabla \left( \nabla F \right) \right] =$$

$$= C_1^2 \left[ \nabla \left( \vec{b}^a \frac{\partial F}{\partial q^a} \right) \right] =$$

$$= C_1^2 \left[ \nabla \left( \vec{b}_b \, g^{ab} \frac{\partial F}{\partial q^a} \right) \right] =$$

$$= C_1^2 \left[ \vec{b}^c \frac{\partial}{\partial q^c} \left( \vec{b}_b \, g^{ab} \frac{\partial F}{\partial q^a} \right) \right] =$$

$$= C_1^2 \left\{ \vec{b}^c \left[ \vec{b}_b \frac{\partial}{\partial q^c} \left( g^{ab} \frac{\partial F}{\partial q^a} \right) + g^{ab} \frac{\partial F}{\partial q^a} \frac{\partial \vec{b}_b}{\partial q^c} \right] \right\} =$$

$$= C_1^2 \left\{ \vec{b}^c \left[ \vec{b}_b \frac{\partial}{\partial q^c} \left( g^{ab} \frac{\partial F}{\partial q^a} \right) + g^{ab} \frac{\partial F}{\partial q^a} \Gamma_{bc}^d \, \vec{b}_d \right] \right\} =$$

$$= C_1^2 \left\{ \vec{b}^c \vec{b}_b \left[ \frac{\partial}{\partial q^c} \left( g^{ab} \frac{\partial F}{\partial q^a} \right) + g^{ad} \Gamma_{cd}^b \frac{\partial F}{\partial q^a} \right] \right\} =$$

$$= \frac{\partial}{\partial q^b} \left( g^{ab} \frac{\partial F}{\partial q^a} \right) + g^{ad} \Gamma_{bd}^b \frac{\partial F}{\partial q^a} \ .$$

---

**Example 9.3.7 (Laplacian of a vector field - with general coordinates $q^a$)**

$$\nabla \cdot \nabla F = C_1^2 \left[ \nabla \left( \nabla F \right) \right] =$$

$$= C_1^2 \left\{ \nabla \left[ \vec{b}^a \, \vec{b}_b \left( \frac{\partial F^b}{\partial q^a} + \Gamma_{ac}^b F^c \right) \right] \right\} =$$

$$= C_1^2 \left\{ \nabla \left[ \vec{b}_c \, \vec{b}_b \, g^{ac} \left( \frac{\partial F^b}{\partial q^a} + \Gamma_{ac}^b F^c \right) \right] \right\} =$$

$$= C_1^2 \left\{ \nabla \cdot \left( (\nabla F)^{cb} \, \vec{b}_c \, \vec{b}_b \right) \right\} =$$

$$= C_1^2 \left\{ \vec{b}^a \, \vec{b}_c \, \vec{b}_b \left[ \frac{\partial (\nabla F)^{cb}}{\partial q^a} + \Gamma_{ad}^c (\nabla F)^{db} + \Gamma_{ad}^b (\nabla F)^{cd} \right] \right\} =$$

$$= \vec{b}_b \left[ \frac{\partial (\nabla F)^{ab}}{\partial q^a} + \Gamma_{ad}^a (\nabla F)^{db} + \Gamma_{ad}^b (\nabla F)^{ad} \right] = \ .$$

---

## 9.3.5 Curl

# 9.4 Integrals in $E^d$, $d \leq 3$

## 9.4.1 Line integrals

### Density

Integrals

$$\int_{\vec{r} \in \gamma} F(\vec{r})$$

represent the summation of contributions $F(\vec{r})$ over elementary segments of path $\gamma$, whose dimension is $|d\vec{r}|$, i.e. implicitly means

$$\int_{\vec{r} \in \gamma} F(\vec{r}) = \int_{\vec{r} \in \gamma} F(\vec{r}) \, |d\vec{r}| \ .$$

Given a regular parametrization of the curve $\vec{r}(q^1)$ (with increasing $q^1$ so that $|dq^1| = dq^1$), and the differential $d\vec{r} = \vec{r}'(q^1) \, dq^1$, the integral can be written as an integral in the parameter $q^1$

$$\int_{q=q_a^1}^{q_b^1} F(\vec{r}(q^1)) \, |\vec{r}'(q^1)| \, dq^1 \ ,$$

with $\vec{r}(q_a^1)$, $\vec{r}(q_b^1)$ the extreme points of path $\gamma$.

### Work

Integrals

$$\int_{\vec{r} \in \gamma} F(\vec{r}) \cdot \hat{t}(\vec{r})$$

implicitly mean

$$\int_{\vec{r} \in \gamma} F(\vec{r}) \cdot \hat{t}(\vec{r}) = \int_{\vec{r} \in \gamma} F(\vec{r}) \cdot \hat{t}(\vec{r}) |d\vec{r}| = \int_{\vec{r} \in \gamma} F(\vec{r}) \cdot d\vec{r} \ ,$$

as $\hat{t} = \frac{d\vec{r}}{|d\vec{r}|}$. Given a regular parametrization of the curve $\vec{r}(q^1)$ (with increasing $q^1$ so that $|dq^1| = dq^1$), and the differential $d\vec{r} = \vec{r}'(q^1) \, dq^1$, the integral can be written as an integral in the parameter $q^1$

$$\int_{q^1=q_a^1}^{q_b^1} F(\vec{r}(q^1)) \cdot \vec{r}'(q^1) \, dq^1$$

## 9.4.2 Surface integrals

Given two coordinates $q^1$, $q^2$ describing a surface, $\vec{r}(q^1, q^2)$ the elementary surface with unit normal reads

$$\hat{n} \, dS = d\vec{r}_1 \times d\vec{r}_2 = \frac{\partial \vec{r}}{\partial q^1} \times \frac{\partial \vec{r}}{\partial q^2} \, dq^1 \, dq^2 \ ,$$

and the elementary surface thus reads

$$|dS| = |\hat{n} dS| = \left| \frac{\partial \vec{r}}{\partial q^1} \times \frac{\partial \vec{r}}{\partial q^2} \, dq^1 \, dq^2 \right|$$

**Density**

Integrals

$$\int_{\vec{r}\in S} F(\vec{r})$$

implicitly mean

$$\int_{\vec{r}\in S} F(\vec{r}) = \int_{\vec{r}\in S} F(\vec{r})|dS| .$$

Given regular parametrization of the surface, $\vec{r}(q^1, q^2)$, $(q^1, q^2) \in Q^{12}$, the integral can be written as the multidimensional integral in coordinates $q^1$, $q^2$,

$$\int_{\vec{r}\in S} F(\vec{r}) = \int_{(q^1,q^2)\in Q^{12}} F(\vec{r}(q^1, q^2))\left|\frac{\partial \vec{r}}{\partial q^1} \times \frac{\partial \vec{r}}{\partial q^2}\, dq^1\, dq^2\right|$$

**Flux**

Integrals

$$\int_{\vec{r}\in S} \hat{n}(\vec{r}) \cdot F(\vec{r})$$

implicitly mean

$$\int_{\vec{r}\in S} \hat{n}(\vec{r}) \cdot F(\vec{r}) = \int_{\vec{r}\in S} \hat{n}(\vec{r}) \cdot F(\vec{r})|dS|$$

Given regular parametrization of the surface, $\vec{r}(q^1, q^2)$, $(q^1, q^2) \in Q^{12}$, the integral can be written as the multidimensional integral in coordinates $q^1$, $q^2$,

$$\int_{\vec{r}\in S} \hat{n}(\vec{r}) \cdot F(\vec{r}) = \int_{(q^1,q^2)\in Q^{12}} \frac{\partial \vec{r}}{\partial q^1} \times \frac{\partial \vec{r}}{\partial q^2} \cdot F(\vec{r}(q^1, q^2))\, dq^1\, dq^2$$

### 9.4.3 Volume

$$dV = \frac{\partial \vec{r}}{\partial q^1} \cdot \frac{\partial \vec{r}}{\partial q^2} \times \frac{\partial \vec{r}}{\partial q^3}\, dq^1\, dq^2\, dq^3 .$$

**Density**

Integrals

$$\int_{\vec{r}\in V} F(\vec{r})$$

implicitly mean

$$\int_{\vec{r}\in V} F(\vec{r}) = \int_{\vec{r}\in V} F(\vec{r})\, |dV| .$$

Given regular parametrization of the volume, $\vec{r}(q^1, q^2, q^3)$, $(q^1, q^2, q^3) \in Q$, the integral can be written as the multidimensional integral in coordinates $q^1$, $q^2$, $q^3$,

$$\int_{\vec{r}\in V} F(\vec{r})|dV| = \int_{(q^1,q^2,q^3)\in Q} F(\vec{r}(q^1, q^2, q^3))\left|\frac{\partial \vec{r}}{\partial q^1} \cdot \frac{\partial \vec{r}}{\partial q^2} \times \frac{\partial \vec{r}}{\partial q^3}\, dq^1\, dq^2\, dq^3\right| .$$

## 9.4.4 Theorems

### Two useful lemmas

The next lemma forms the foundation of the well-known *divergence theorem* and *gradient theorem*:
the proof of these two theorems is based on a straightforward repeated use of this lemma.
Given how simple this lemma is and how frequently it is applied in writing balances and, more generally, in integration
by parts, it is very convenient to remember this simple result.

---

**Theorem 9.4.1 (Lemma 1.)**

Under the assumptions of *Green's lemma in the plane*,

$$\int_V \frac{\partial A}{\partial x_i} = \oint_S A n_i$$

---

### Proof

The reasoning closely follows the one used for the proof of *Green's lemma in the plane*.
For $\partial A / \partial z$:

$$\int_V \frac{\partial A}{\partial z} = \int_R \int_{z=f_1(x,y)}^{z=f_2(x,y)} \frac{\partial A}{\partial z} dz dx dy =$$

$$= \int_R [A(x, y, f_2(x,y)) - A(x, y, f_1(x,y))] dx dy$$

The most complex step is transitioning from the integral over $(x, y) \in R$ to the surface integral over $S$, the boundary of
volume $V$: the infinitesimal area element $dR$ in the xy-plane is equal to $dR = dx dy$; the drawing and the proof refer
to a *simple* volume (as in the case of Green's lemma in the plane, the results can be generalized to domains of arbitrary
shape).
It is possible to divide the surface $S$ into two "halves" $S^+ : z = f_2(x,y)$ and $S^- : z = f_1(x, y)$ such that
$S^+ \cup S^- = S$, and the outward normal has positive and negative z-components respectively ($S^+ : \hat{\mathbf{n}} \cdot \hat{\mathbf{z}} > 0, S^- : \hat{\mathbf{n}} \cdot \hat{\mathbf{z}} < 0$).
The surface element $dR$ is also the projection of the surface element $dS$ onto the xy-plane: in general, $dS$ will not be
parallel to the xy-plane and thus will be larger than $dR$. It's not difficult to show that:

$$dx dy = dR = \begin{cases} dS \hat{\mathbf{z}} \cdot \hat{\mathbf{n}} & \text{on } S^+ \\ -dS \hat{\mathbf{z}} \cdot \hat{\mathbf{n}} & \text{on } S^- \end{cases}$$

We can now continue the proof:

$$\int_R [A(x, y, f_2(x,y)) - A(x, y, f_1(x,y))] dx dy =$$

$$= \int_{S^+} A \hat{\mathbf{n}} \cdot \hat{\mathbf{z}} dS + \int_{S^-} A \hat{\mathbf{n}} \cdot \hat{\mathbf{z}} dS =$$

$$= \oint_S A \hat{\mathbf{z}} \cdot \hat{\mathbf{n}} dS =$$

$$= \oint_S A n_z dS$$

Just as the previous lemma forms the basis for the proof of the *gradient theorem* and the *divergence theorem*,
the following lemma forms the basis for the proof of the *curl theorem*.

---

**Theorem 9.4.2 (Lemma 2.)**

Under the assumptions of *Green's lemma in the plane*,

$$\int_S [\mathbf{\nabla} \times (A\hat{\mathbf{e}}_\mathbf{i})] \cdot \hat{\mathbf{n}} = \oint_\gamma A\, dx_i$$

**Proof**

For $A\hat{\mathbf{e}}_\mathbf{x}$, we have $\nabla \times (A\hat{\mathbf{e}}_\mathbf{x}) = \partial A/\partial z\hat{\mathbf{e}}_\mathbf{y} - \partial A/\partial y\hat{\mathbf{e}}_\mathbf{z}$. The surface $S$ is written in parametric form as: $\mathbf{r} = x\hat{\mathbf{e}}_\mathbf{x} + y\hat{\mathbf{e}}_\mathbf{y} + z(x,y)\hat{\mathbf{e}}_\mathbf{z}$. The vector $\partial\mathbf{r}/\partial y = \hat{\mathbf{e}}_\mathbf{y} + \partial z/\partial y\hat{\mathbf{e}}_\mathbf{z}$ is tangent to the surface $S$ and hence perpendicular to the normal $\hat{\mathbf{n}}$:

$$0 = \hat{\mathbf{n}} \cdot \left( \hat{\mathbf{e}}_\mathbf{y} + \frac{\partial z}{\partial y}\hat{\mathbf{e}}_\mathbf{z} \right)$$

Now writing $[\mathbf{\nabla} \times (A\hat{\mathbf{e}}_\mathbf{x})] \cdot \hat{\mathbf{n}}$:

$$[\mathbf{\nabla} \times (A\hat{\mathbf{e}}_\mathbf{x})] \cdot \hat{\mathbf{n}} = \frac{\partial A}{\partial z}\hat{\mathbf{e}}_\mathbf{y} \cdot \hat{\mathbf{n}} - \frac{\partial A}{\partial y}\hat{\mathbf{e}}_\mathbf{z} \cdot \hat{\mathbf{n}} = -\left[ \frac{\partial A}{\partial z}\frac{\partial z}{\partial y} + \frac{\partial A}{\partial y} \right]\hat{\mathbf{e}}_\mathbf{z} \cdot \hat{\mathbf{n}}$$

Recognizing that $\frac{\partial A(x,y,z(x,y))}{\partial y} = \frac{\partial A}{\partial z}\frac{\partial z}{\partial y} + \frac{\partial A}{\partial y}$, we can write:

$$\int_S [\mathbf{\nabla} \times (A\hat{\mathbf{e}}_\mathbf{x})] \cdot \hat{\mathbf{n}} = -\int_S \frac{\partial A}{\partial y} \underbrace{\hat{\mathbf{e}}_\mathbf{z} \cdot \hat{\mathbf{n}}\,dS}_{dR=dxdy} = -\int_R \frac{\partial A}{\partial y}\,dx\,dy = \int_\gamma A\,dx$$

**Gradient theorem**

…**todo** *assumptions*…

$$\int_V \nabla f = \oint_{\partial V} f\hat{n}$$

**Proof for simple domains** $V$

This result immediately follows from Lemma 1 *Theorem 9.4.1*

$$\oint_{\partial V} f\hat{n} = \hat{x}_i \oint_{\partial V} fn_i = \hat{x}_i \int_V \partial_i f = \int_V \hat{x}_i \partial_i f = \int_V \nabla\vec{f}\,,$$

having (1) exploited the freedom to put unit vectors of the Cartesian basis inside the integrals, as they're unifrom - constant in space -, and (2) recognized the expression of the *gradient of a scalar field expressed using Cartesian coordinates*, as shown in *Example 9.5.1*.

**Divergence theorem**

…**todo** *assumptions*…

$$\int_V \nabla \cdot \vec{f} = \oint_{\partial V} \vec{f} \cdot \hat{n}$$

### Proof

This result immediately follows from Lemma 1 *Theorem 9.4.1*

$$\oint_{\partial V} \vec{f} \cdot \hat{n} = \oint_{\partial V} f_i n_i = \int_V \partial_i f_i = \int_V \nabla \cdot \vec{f} \,,$$

having (1) exploited the freedom to put unit vectors of the Cartesian basis inside the integrals, as they're unifrom - constant in space -, and (2) recognized the expression of the *divergence of a vector field expressed using Cartesian coordinates*, as shown in *Example 9.5.4*.

### Curl theorem

…**todo** *assumptions*…

$$\int_S \left[ \nabla \times \vec{f} \right] \cdot \hat{n} = \oint_{\partial S} \vec{f} \cdot \hat{t}$$

### Proof

This proof seamlessly follows from Lemma *Theorem 9.4.2*, applied to all the Cartesian contributions of the vector field

$$\vec{f} = f_x \hat{x} + f_y \hat{y} + f_z \hat{z} \,,$$

as

$$\int_S \left( \nabla \times \vec{f} \right) \cdot \hat{n} = \int_S \left( \nabla \times f_i \hat{x}_i \right) \cdot \hat{n} = \oint_{\partial S} f_i t_i = \oint_{\partial S} \vec{f} \cdot \hat{t} \,.$$

## 9.5 Tensor Calculus in Euclidean Spaces - Cartesian coordinates in $E^3$

Using Cartesian coordinates $(q^1, q^2, q^3) = (r, \theta, z)$ and Cartesian base vectors (uniform in space, so that their derivatives are zero), a point in Euclidean vector space $E^3$ can be represented as

$$\vec{r} = x \, \hat{x} + y \, \hat{y} + z \, \hat{z} \,.$$

### 9.5.1 Natural basis, reciprocal basis, metric tensor, and Christoffel symbols

Cartesian coordinates in Euclidean spaces are a very special coordinate system, with reciprocal basis everywhere coinciding with natural basis, with uniform basis in space (zero second-order derivative of space w.r.t. coordinates, and thus zero first order derivative of base vectors, and thus identically zero Christoffel symbols), and components of the metric tensor equal to the identity matrix

$$\begin{cases} \vec{b}_1 = \vec{b}^1 = \hat{x} \\ \vec{b}_2 = \vec{b}^2 = \hat{y} \\ \vec{b}_3 = \vec{b}^3 = \hat{z} \end{cases}$$

$$[g_{ab}] = [g^{ab}] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\Gamma_{ab}^c = 0 \quad , \quad \forall a, b, c = 1 : 3 \,.$$

## 9.5.2 Differential operators

### Gradient

---

**Example 9.5.1 (Gradient of a scalar field)**

$$\nabla F = \hat{x}\,\partial_x F + \hat{y}\,\partial_y F + \hat{z}\,\partial_z F_z$$

---

**Example 9.5.2 (Gradient of a vector field)**

$$\begin{aligned}
\nabla F = \nabla(F_x\hat{x} + F_y\hat{y} + F_z\hat{z}) &= \\
&= \hat{x}\otimes\hat{x}\,\partial_x F_x + \hat{x}\otimes\hat{y}\,\partial_x F_y + \hat{x}\otimes\hat{z}\,\partial_x F_z + \\
&\quad + \hat{y}\otimes\hat{x}\,\partial_y F_x + \hat{y}\otimes\hat{y}\,\partial_y F_y + \hat{y}\otimes\hat{z}\,\partial_y F_z + \\
&\quad + \hat{z}\otimes\hat{x}\,\partial_z F_x + \hat{z}\otimes\hat{y}\,\partial_z F_y + \hat{z}\otimes\hat{z}\,\partial_z F_z +
\end{aligned}$$

---

**Example 9.5.3 (Gradient of a $2^{nd}$-order tensor field)**

---

### Directional derivative

### Divergence

---

**Example 9.5.4 (Divergence of a vector field)**

$$\begin{aligned}
\nabla \cdot F = \nabla \cdot \left(F_x\,\hat{x} + F_y\,\hat{y} + F_z\,\hat{z}\right) &= \\
&= \partial_x F_x + \partial_y F_y + \partial_z F_z \;.
\end{aligned}$$

---

**Example 9.5.5 (Divergence of a $2^{nd}$-order tensor field)**

$$\begin{aligned}
\nabla \cdot F = \nabla \cdot (F_{ab}\vec{e}_a \otimes \vec{e}_b) &= \\
&= \vec{e}_c \frac{\partial F_{ab}}{\partial q^a} = \\
&= \hat{x}\left[\partial_x F_{xx} + \partial_y F_{yx} + \partial_z F_{zx}\right] + \\
&\quad + \hat{y}\left[\partial_x F_{xy} + \partial_y F_{yy} + \partial_z F_{zy}\right] + \\
&\quad + \hat{z}\left[\partial_x F_{xz} + \partial_y F_{yz} + \partial_z F_{zz}\right] \;.
\end{aligned}$$

---

**Laplacian**

---

**Example 9.5.6 (Laplacian of a scalar field)**

$$\nabla^2 F = \partial_{xx} F + \partial_{yy} F + \partial_{zz} F$$

---

**Example 9.5.7 (Laplacian of a vector field)**

---

# 9.6 Tensor Calculus in Euclidean Spaces - cylindrical coordinates in $E^3$

## 9.6.1 Cylindrical coordiantes and cylindrical coordinates

Using cylindrical coordinates $(q^1, q^2, q^3) = (r, \theta, z)$ and cylindrical base vectors (uniform in space, so that their derivatives are zero), a point in Euclidean vector space $E^3$ can be represented as

$$\vec{r} = r \cos\theta\, \hat{x} + r \sin\theta\, \hat{y} + z\, \hat{z} \ .$$

## 9.6.2 Natural basis, reciprocal basis, metric tensor, and Christoffel symbols

### Natural basis

Natural basis reads

$$\begin{cases} \vec{b}_1 = \dfrac{\partial \vec{r}}{\partial q^1} = \dfrac{\partial \vec{r}}{\partial r} = \cos\theta\, \hat{x} + \sin\theta\, \hat{y} \\[2mm] \vec{b}_2 = \dfrac{\partial \vec{r}}{\partial q^2} = \dfrac{\partial \vec{r}}{\partial \theta} = -r \sin\theta\, \hat{x} + r \cos\theta\, \hat{y} \\[2mm] \vec{b}_3 = \dfrac{\partial \vec{r}}{\partial q^3} = \dfrac{\partial \vec{r}}{\partial z} = \hat{z} \end{cases}$$

### Metric tensor

Covariant components of metric tensors,

$$g_{ab} = \vec{b}_a \cdot \vec{b}_b \ ,$$

can be collected in the diagonal matrix

$$[g_{ab}] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & r^2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \ ,$$

while its contra-variant components can be collected in the inverse matrix (easy to compute, since $[g_{ab}]$ is diagonal),

$$[g^{ab}] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{r^2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \ .$$

## Reciprocal basis

Reciprocal basis is readily evaluated using $\vec{b}^a = g^{ab}\,\vec{b}_b$,

$$\begin{cases} \vec{b}^1 = \cos\theta\,\hat{x} + \sin\theta\,\hat{y} \\ \vec{b}^2 = -\dfrac{1}{r}\sin\theta\,\hat{x} + \dfrac{1}{r}\cos\theta\,\hat{y} \\ \vec{b}^3 = \hat{z} \end{cases}$$

## Physical basis

Since metric tensor is diagonal, the cylindrical coordinate system is orthogonal, and its natural and reciprocal basis are orthogonal. A unit orthogonal basis, usually named **physical basis** with unit vector with no physical dimension, is evalated by normalization process,

$$\begin{cases} \hat{r} = \hat{b}_1 = \dfrac{\vec{b}_1}{g_{11}} = \dfrac{\vec{b}^1}{g^{11}} = \cos\theta\,\hat{x} + \sin\theta\,\hat{y} \\ \hat{\theta} = \hat{b}_2 = \dfrac{\vec{b}_2}{g_{22}} = \dfrac{\vec{b}^2}{g^{22}} = -\sin\theta\,\hat{x} + \cos\theta\,\hat{y} \\ \hat{z} = \hat{b}_3 = \dfrac{\vec{b}_3}{g_{33}} = \dfrac{\vec{b}^3}{g^{33}} = \hat{z}\;. \end{cases}$$

## Derivatives of natural basis and Christoffel symbols

Derivatives of the natural basis read

$$\frac{\partial \vec{b}_1}{\partial q^1} = \vec{0}$$

$$\frac{\partial \vec{b}_2}{\partial q^2} = -r\cos\theta\,\hat{x} - r\sin\theta\,\hat{y} = -q^1\,\vec{b}_1$$

$$\frac{\partial \vec{b}_3}{\partial q^3} = \vec{0}$$

$$\frac{\partial \vec{b}_2}{\partial q^1} = \frac{\partial \vec{b}_1}{\partial q^2} = -\sin\theta\,\hat{x} + \cos\theta\hat{y} = \frac{1}{q^1}\,\vec{b}_2$$

$$\frac{\partial \vec{b}_3}{\partial q^1} = \frac{\partial \vec{b}_1}{\partial q^3} = \vec{0}$$

$$\frac{\partial \vec{b}_3}{\partial q^2} = \frac{\partial \vec{b}_2}{\partial q^3} = \vec{0}$$

so that non-zero Christoffel symbols of a cylindrical coordinate system are

$$\Gamma^2_{12} = \Gamma^2_{21} = \frac{1}{q^1}$$

$$\Gamma^1_{22} = -q^1\;.$$

## 9.6.3 Differential operators

### Gradient

**Example 9.6.1 (Gradient of a scalar field)**

$$\nabla F = \vec{b}^a \frac{\partial F}{\partial q^a} =$$

$$= \hat{b}_a \, g^{aa} \, \frac{\partial F}{\partial q^a} =$$

$$= \hat{r} \frac{\partial F}{\partial r} + \hat{\theta} \frac{1}{r} \frac{\partial F}{\partial \theta} + \hat{z} \frac{\partial F}{\partial z} \, .$$

**Example 9.6.2 (Gradient of a vector field)**

$$\nabla F = \vec{b}^a \otimes \vec{b}_b \left[ \frac{\partial F^b}{\partial q^a} + \Gamma^b_{ac} \, F^c \right] =$$

$$= \cdots =$$

$$
\begin{aligned}
&= \vec{b}^1 \otimes \vec{b}_1 \, \partial_1 F^1 && + \vec{b}^1 \otimes \vec{b}_2 \, [\partial_1 F^2 + \Gamma^2_{12} F^2] && + \vec{b}^1 \otimes \vec{b}_3 \, \partial_1 F^3 \\
&+ \vec{b}^2 \otimes \vec{b}_1 \, [\partial_2 F^1 + \Gamma^1_{22} F^2] && + \vec{b}^2 \otimes \vec{b}_2 \, [\partial_2 F^2 + \Gamma^2_{21} F^1] && + \vec{b}^2 \otimes \vec{b}_3 \, \partial_2 F^3 \\
&+ \vec{b}^3 \otimes \vec{b}_1 \, \partial_3 F^1 && + \vec{b}^3 \otimes \vec{b}_2 \, \partial_3 F^2 && + \vec{b}^3 \otimes \vec{b}_3 \, \partial_3 F^3 \\
&= \hat{r} \otimes \hat{r} \, \partial_r F_r && + \hat{r} \otimes \hat{\theta} \, \frac{1}{r} \, [\partial_r (r F_\theta) + F_\theta] && + \hat{r} \otimes \hat{z} \, \partial_r F_z \\
&+ \hat{\theta} \otimes \hat{r} \, \frac{1}{r} \, [\partial_\theta F_r - r F_\theta] && + \hat{\theta} \otimes \hat{\theta} \left[ \partial_\theta \left( \frac{F_\theta}{r} \right) + \frac{F_r}{r} \right] && + \hat{\theta} \otimes \hat{z} \, \frac{1}{r} \partial_\theta F_z \\
&+ \hat{z} \otimes \hat{r} \, \partial_z F_x && + \hat{z} \otimes \hat{\theta} \, \frac{1}{r} \partial_\theta F_y && + \hat{z} \otimes \hat{z} \, \partial_z F_z \, .
\end{aligned}
$$

**Example 9.6.3 (Gradient of a $2^{nd}$-order tensor field)**

### Directional derivative

### Divergence

**Example 9.6.4 (Divergence of a vector field)**

$$\nabla \cdot \vec{F} = \frac{\partial F^a}{\partial q^a} + \Gamma^a_{ac} \, F^c =$$

$$= \frac{\partial F_r}{\partial r} + \frac{\partial}{\partial \theta} \left( \frac{F_\theta}{r} \right) + \frac{F_\theta}{r} + \frac{\partial F_z}{\partial z} \, .$$

### Divergence of a $2^{nd}$-order tensor field

Using the general formula of the divergence of a $2^{nd}$-order tensor field (see *Divergence*)

$$\nabla \cdot \left( F^{bc} \vec{b}_b \otimes \vec{b}_c \right) = C_1^2 \left( \nabla F \right) =$$

$$= \vec{b}_c \left[ \frac{\partial F^{ac}}{\partial q^a} + \Gamma^a_{ad} F^{dc} + \Gamma^c_{ad} F^{ad} \right]$$

the contravariant components in the natural basis induced by cylindrical coordinates of the divergence of a second order tensor reads

$$\nabla \cdot \mathbb{F} =$$

$$= \vec{b}_1 \left[ \frac{\partial F^{11}}{\partial q^1} + \frac{\partial F^{21}}{\partial q^2} + \frac{\partial F^{31}}{\partial q^3} + \Gamma^2_{21} F^{11} + \Gamma^1_{22} F^{22} \right] +$$

$$+ \vec{b}_2 \left[ \frac{\partial F^{12}}{\partial q^1} + \frac{\partial F^{22}}{\partial q^2} + \frac{\partial F^{32}}{\partial q^3} + \Gamma^2_{21} F^{12} + \Gamma^2_{21} F^{12} + \Gamma^2_{21} F^{21} \right] +$$

$$+ \vec{b}_3 \left[ \frac{\partial F^{13}}{\partial q^1} + \frac{\partial F^{23}}{\partial q^2} + \frac{\partial F^{33}}{\partial q^3} + \Gamma^2_{21} F^{13} \right] =$$

$$= \vec{b}_1 \left[ \frac{\partial F^{11}}{\partial q^1} + \frac{\partial F^{21}}{\partial q^2} + \frac{\partial F^{31}}{\partial q^3} + \frac{1}{q^1} F^{11} - q^1 F^{22} \right] +$$

$$+ \vec{b}_2 \left[ \frac{\partial F^{12}}{\partial q^1} + \frac{\partial F^{22}}{\partial q^2} + \frac{\partial F^{32}}{\partial q^3} + \frac{2}{q^1} F^{12} + \frac{1}{q^1} F^{21} \right] +$$

$$+ \vec{b}_3 \left[ \frac{\partial F^{13}}{\partial q^1} + \frac{\partial F^{23}}{\partial q^2} + \frac{\partial F^{33}}{\partial q^3} + \frac{1}{q^1} F^{13} \right] .$$

Next, it's easy to exploit the definition of the coordinates $(q^1, q^2, q^3) = (r, \theta, z)$ and the relation between natural and physical basis and components to get

$$\nabla \cdot \mathbb{F} =$$

$$= \hat{r} \left[ \frac{\partial \left( F^{rr} \right)}{\partial r} + \frac{\partial \left( \frac{1}{r} F^{\theta r} \right)}{\partial \theta} + \frac{\partial F^{zr}}{\partial z} + \frac{1}{r} F^{rr} - r \frac{1}{r^2} F^{\theta\theta} \right] +$$

$$+ r\hat{\theta} \left[ \frac{\partial \left( \frac{1}{r} F^{r\theta} \right)}{\partial r} + \frac{1}{r} \frac{\partial \left( \frac{1}{r} F^{\theta\theta} \right)}{\partial \theta} + \frac{\partial F^{z\theta}}{\partial z} + \frac{2}{r} \frac{1}{r} F^{r\theta} + \frac{1}{r} \frac{1}{r} F^{\theta r} \right] +$$

$$+ \hat{z} \left[ \frac{\partial F^{rz}}{\partial r} + \frac{\partial \left( \frac{1}{r} F^{\theta z} \right)}{\partial \theta} + \frac{\partial F^{zz}}{\partial z} + \frac{1}{r} F^{rz} \right] .$$

or, after few algebraic manipulations,

$$\nabla \cdot \mathbb{F} =$$

$$= \hat{r} \left[ \frac{\partial F^{rr}}{\partial r} + \frac{1}{r} \frac{\partial F^{\theta r}}{\partial \theta} + \frac{\partial F^{zr}}{\partial z} + \frac{1}{r} F^{rr} - \frac{1}{r} F^{\theta\theta} \right] +$$

$$+ \hat{\theta} \left[ \frac{\partial F^{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial F^{\theta\theta}}{\partial \theta} + \frac{\partial F^{z\theta}}{\partial z} + \frac{1}{r} F^{r\theta} + \frac{1}{r} F^{\theta r} \right] +$$

$$+ \hat{z} \left[ \frac{\partial F^{rz}}{\partial r} + \frac{1}{r} \frac{\partial F^{\theta z}}{\partial \theta} + \frac{\partial F^{zz}}{\partial z} + \frac{1}{r} F^{rz} \right] .$$

**Laplacian**

---

**Example 9.6.5 (Laplacian of a scalar field)**

---

**Example 9.6.6 (Laplacian of a vector field)**

---

# 9.7 Tensor Calculus in Euclidean Spaces - Spehrical coordinates in $E^3$

Using spherical coordinates $(q^1, q^2, q^3) = (r, \phi, \theta)$ and spherical base vectors (uniform in space, so that their derivatives are zero), a point in Euclidean vector space $E^3$ can be represented as

$$\vec{r} = r \cos \theta \sin \phi \, \hat{x} + r \sin \theta \sin \phi \, \hat{y} + r \cos \phi \, \hat{z} \, .$$

## 9.7.1 Natural basis, reciprocal basis, metric tensor, and Christoffel symbols

## 9.7.2 Differential operators

**Gradient**

---

**Example 9.7.1 (Gradient of a scalar field)**

---

**Example 9.7.2 (Gradient of a vector field)**

---

**Example 9.7.3 (Gradient of a $2^{nd}$-order tensor field)**

---

**Directional derivative**

**Divergence**

---

**Example 9.7.4 (Divergence of a vector field)**

---

**Example 9.7.5 (Divergence of a $2^{nd}$-order tensor field)**

---

## Laplacian

**Example 9.7.6 (Laplacian of a scalar field)**

**Example 9.7.7 (Laplacian of a vector field)**

# TENSOR INVARIANTS

## 10.1 Rank-$2$ tensors

- Characteristic polynomial (*this definition needs that determinant exists*)

$$0 = \det\left(\mathbf{A} - s\mathbf{I}\right) \ .$$

*Using a unit normal basis* (**todo** generalization required?)

$$0 = \begin{vmatrix} A_{11} - s & A_{12} & A_{13} \\ A_{21} & A_{22} - s & A_{23} \\ A_{31} & A_{32} & A_{33} - s \end{vmatrix} = -s^3 + s^2 I_1 - s I_2 + I_3 \ ,$$

with

$$I_1 = \text{tr}(\mathbf{A}) = A_{11} + A_{22} + A_{33}$$

$$I_2 = \frac{1}{2}\left(\text{tr}(\mathbf{A})^2 - \text{tr}(\mathbf{A}^2)\right) = \cdots = A_{11}A_{22} + A_{11}A_{33} + A_{22}A_{33} - A_{31}A_{13} - A_{12}A_{21} - A_{32}A_{23}$$

$$I_3 = \det(\mathbf{A}) = ...$$

The coefficients of the characteristic polynomial are invariant under transformation of coordinates (**todo** here only referring to Cartesian basis, and thus orthogonal transformations?)

- **Trace**…

- Determinant…

# UNITARY AND ROTATION TENSORS

Some notes in the introduction to classical mechanics: Rotations: Tensor formalism for rotations.

**todo** Move the mathematical treatment here?

**todo** Discuss vector operations that are invariant under 3-dimensional rotation and other unitary tensor applications (i.e. reflections). These properties are useful to prove the general expression of *isotropic tensors*.

## 11.1 Invariant operations under rotations and unitary transformations

Proof that/if:

- the only independent invariant operations producing a scalar with **two** vectors $\mathbf{u}$, $\mathbf{v}$ are

$$\mathbf{u} \cdot \mathbf{u} = |\mathbf{u}|^2$$
$$\mathbf{v} \cdot \mathbf{v} = |\mathbf{v}|^2$$
$$\mathbf{u} \cdot \mathbf{v}$$

- the only independent invariant operations producing a scalar with **three** vectors $\mathbf{u}$, $\mathbf{v}$, $\mathbf{w}$ are

$$\mathbf{u} \times \mathbf{v} \cdot \mathbf{w}$$

- the only independent invariant operations producing a scalar with **four** vectors $\left\{\mathbf{u}^{(i)}\right\}_{i=1:4}$ are

$$\mathbf{u}^{(k_1)} \cdot \mathbf{u}^{(k_2)} \, \mathbf{u}^{(k_3)} \cdot \mathbf{u}^{(k_4)} \ ,$$

with every index $k_j$ independently ranging from $1$ to $4$.

# ISOTROPIC TENSORS

This section provides the definition of isotropic tensors, the general forms of low-rank tensors (from rank-2 to rank-4 tensors) that may be useful in some fields of science (e.g. in continuum mechanics, in the constitutive law relating stress tensor with strain tensor in isotropic linear elastic media or stress tensor with strain velocity tensor in Newtonian fluids).

---

**Definition 12.1 (Isotropic tensor)**

An isotropic tensor is …

---

The general form of the low-rank isotropic tensors in three-dimensional spaces, $\dim(\mathcal{V}) = 3$ is proved using the definition of a rank-$r$ tensor as a multilinear map acting on $r$ vectors, and exploiting the invariance of some vector operations under rotations.

## 12.1 Rank-$2$ isotropic tensors

The most general rank-2 isotropic tensor is a scalar multiple of the identity tensor, $\mathbf{A} = a\mathbf{I}$.

### Proof

The action of a rank-2 isotropic tensor in $\mathcal{V}$ on every pair of vectors $\mathbf{u}$, $\mathbf{v}$ can be a function of $\mathbf{u} \cdot \mathbf{v}$ only, i.e. the only vector operation:

- producing a scalar

- invariant under rotations

- involving the two vectors $\mathbf{u}$, $\mathbf{v}$ (i.e. the tensor is a multilinear map, $\mathbf{A}(\mathbf{u}, \mathbf{v})$)

Now, using coordinates, the general expression of the isotropic tensor is retrieved comparing the expression built under the isotropic assumption,

$$\mathbf{A}(\mathbf{u}, \mathbf{v}) = a\mathbf{u} \cdot \mathbf{v} =$$
$$= a\left(u^i \mathbf{b}_i\right) \cdot \left(v^j \mathbf{b}_j\right) =$$
$$= au^i v^j g_{ij}$$

with the general expression of the action of the stress tensor on the vectors $\mathbf{u}$, $\mathbf{v}$

$$\mathbf{A}(\mathbf{u}, \mathbf{v}) = A_{ij}\mathbf{b}^i\mathbf{b}^j\left(u^k\mathbf{b}_k, v^l\mathbf{b}_l\right) = A_{ij}\delta^i_k\delta^j_l u^k v^l = A_{ij}u^i v^j \ ,$$

so that the *covariant components* read

$$A_{ij} = ag_{ij}$$

---

and the tensor can be written as

$$\mathbf{A} = A_{ij}\mathbf{b}^i\mathbf{b}^j = ag_{ij}\mathbf{b}^i\mathbf{b}^j =$$
$$= a\mathbf{b}^i\mathbf{b}_i = a\delta_i^j\mathbf{b}^i\mathbf{b}_j =$$
$$= ag^{ij}\mathbf{b}_j\mathbf{b}_i = ag^{ij}\mathbf{b}_i\mathbf{b}_j .$$

### Components of the identity tensor

The identity tensor can be defined as a rank-$2$ tensor whose dot product with any arbitrary $\mathbf{v} \in \mathbf{V}$ gives

$$\mathbf{I} \cdot \mathbf{v} = \mathbf{v} .$$

The identity tensor can be written using any vector basis $\{\mathbf{b}_i\}_i$, and its reciprocal $\{\mathbf{b}^j\}$ as

$$\mathbf{I} = \mathbf{b}_i\mathbf{b}^i = \mathbf{b}^i\mathbf{b}_i = g_{ij}\mathbf{b}^i\mathbf{b}^j = g^{ij}\mathbf{b}_i\mathbf{b}_j ,$$

as it's immediate to prove by direct computation

$$\mathbf{b}_i\mathbf{b}^i \cdot \mathbf{v} = \mathbf{b}_i\mathbf{b}^i \cdot (v^k\mathbf{b}_k) = \mathbf{b}_i\delta_k^i v^k = v^i\mathbf{b}_i = \mathbf{v}$$
$$\mathbf{b}^i\mathbf{b}_i \cdot \mathbf{v} = \mathbf{b}^i\mathbf{b}_i \cdot (v_k\mathbf{b}^k) = \mathbf{b}^i\delta_i^k v_k = v_i\mathbf{b}^i = \mathbf{v}$$
$$g_{ij}\mathbf{b}^i\mathbf{b}^j \cdot \mathbf{v} = g_{ij}\mathbf{b}^i\mathbf{b}^j \cdot (v^k\mathbf{b}_k) = \mathbf{b}^i\delta_k^j v^k g_{ij} = \mathbf{b}^i \underbrace{v^k g_{ik}}_{=v_i} = \mathbf{v}$$
$$g^{ij}\mathbf{b}_i\mathbf{b}_j \cdot \mathbf{v} = g^{ij}\mathbf{b}_i\mathbf{b}_j \cdot (v_k\mathbf{b}^k) = \mathbf{b}_i\delta_j^k v_k g^{ij} = \mathbf{b}_i \underbrace{v_k g^{ik}}_{=v^i} = \mathbf{v}$$

## 12.2 Rank-$3$ isotropic tensors

**todo**

- Invariant object: $\mathbf{u} \times \mathbf{v} \cdot \mathbf{w}$

- Retrieve components, to find Levi-Civita symbols

- Discuss invariance under rotation and symmetry (pseudo-tensor)

## 12.3 Rank-$4$ isotropic tensors

**Proof**

The action of a rank-$4$ isotropic tensor in $\mathcal{V}$ on every set of 4 vectors $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x}$ can be a function of the products of scalar products

$$\mathbf{u} \cdot \mathbf{v}\,\mathbf{w} \cdot \mathbf{x}$$
$$\mathbf{u} \cdot \mathbf{w}\,\mathbf{v} \cdot \mathbf{x}$$
$$\mathbf{u} \cdot \mathbf{x}\,\mathbf{v} \cdot \mathbf{w} ,$$

i.e. any possible combinations through scalar products of different vectors

- producing a scalar

- invariant under rotations

- involving all the 4 vectors $\mathbf{u}$, $\mathbf{v}$, $\mathbf{w}$, $\mathbf{x}$ (i.e. the tensor is a multilinear map, $\mathbf{A}(\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x})$)

Now, using coordinates, the general expression of the isotropic tensor is retrieved comparing the expression built under the isotropic assumption,

$$\mathbf{A}(\mathbf{u}, \mathbf{v}) = a\,\mathbf{u} \cdot \mathbf{v}\,\mathbf{w} \cdot \mathbf{x} + b\,\mathbf{u} \cdot \mathbf{w}\,\mathbf{v} \cdot \mathbf{x} + c\,\mathbf{u} \cdot \mathbf{x}\,\mathbf{v} \cdot \mathbf{w} =$$
$$= a\,g_{ij}g_{kl}u^i v^j w^k x^l + b\,g_{ij}g_{kl}u^i w^j v^k x^l + c\,g_{ij}g_{kl}x^i v^j v^k w^l =$$
$$= \left(a\,g_{ij}g_{kl} + b\,g_{ik}g_{jl} + c\,g_{il}g_{jk}\right) u^i v^j w^k x^l \; ,$$

with the general expression of the action of the stress tensor on the vectors $\mathbf{u}$, $\mathbf{v}$

$$\mathbf{A}(\mathbf{u}, \mathbf{v}) = A_{ijkl}\mathbf{b}^i\mathbf{b}^j\mathbf{b}^k\mathbf{b}^l \left(u^m\mathbf{b}_m, v^n\mathbf{b}_n, u^p\mathbf{b}_p, v^q\mathbf{b}_q\right) =$$
$$= A_{ijkl}\delta^i_m\delta^j_n\delta^k_p\delta^l_q u^m v^n w^p x^q = A_{ijkl}u^i v^j w^k x^l \; ,$$

so that the *covariant components* read

$$A_{ijkl} = a\,g_{ij}g_{kl} + b\,g_{ik}g_{jl} + c\,g_{il}g_{jk} \; ,$$

and the tensor can be written as

$$\mathbf{A} = A_{ijkl}\mathbf{b}^i\mathbf{b}^j\mathbf{b}^k\mathbf{b}^l = \left(a\,g_{ij}g_{kl} + b\,g_{ik}g_{jl} + c\,g_{il}g_{jk}\right)\mathbf{b}^i\mathbf{b}^j\mathbf{b}^k\mathbf{b}^l =$$

The most general liner isotropic relation between two rank-2 tensors reads $\mathbf{A}$, $\mathbf{B}$

$$\mathbf{A} = a\operatorname{tr}(\mathbf{B})\,\mathbf{I} + b\,\mathbf{B} + c\,\mathbf{B}^T$$
$$= a\operatorname{tr}(\mathbf{B})\,\mathbf{I} + \tilde{b}\,\mathbf{B}^s + \tilde{c}\,\mathbf{B}^a \; ,$$

being $\mathbf{B}^s = \frac{1}{2}\left(\mathbf{B} + \mathbf{B}^T\right)$, and $\mathbf{B}^a = \frac{1}{2}\left(\mathbf{B} - \mathbf{B}^T\right)$ the symmetric and the anti-symmetric parts of the tensor $\mathbf{B}$ respectively.

### Linear isotropic relation between two rank-$2$ tensors

$$\mathbf{A} = \mathbf{Q} : \mathbf{B}$$
$$\mathbf{Q} : \mathbf{B} = \left(Q_{ijkl}\mathbf{b}^i\mathbf{b}^j\mathbf{b}^k\mathbf{b}^l\right) : \left(B^{mn}\mathbf{b}_m\mathbf{b}_n\right) =$$
$$= Q_{ijkl}B^{kl}\mathbf{b}^i\mathbf{b}^j =$$
$$= \left(a\,g_{ij}g_{kl} + b\,g_{ik}g_{jl} + c\,g_{il}g_{jk}\right)B^{kl}\mathbf{b}^i\mathbf{b}^j =$$
$$= \left(a\,g_{ij}B_l^{\;l} + b\,B_{ij} + c\,B_{ji}\right)\mathbf{b}^i\mathbf{b}^j =$$
$$= \left(a\,g_{ij}B_l^{\;l} + b\,B_{ij} + c\,B_{ji}\right)\mathbf{b}^i\mathbf{b}^j =$$

depends on the three constant $a$, $b$, $c$ of the isotropic tensor.

The most general liner isotropic relation between two symmetric rank-2 tensors reads $\mathbf{A}$, $\mathbf{B}$

$$\mathbf{A} = a\operatorname{tr}(\mathbf{B})\,\mathbf{I} + \tilde{b}\,\mathbf{B} \; ,$$

as the anti-symmetric part of a symmetric tensor is identically zero.

### Linear isotropic relation between two symmetric rank-$2$ tensors

If $B_{ij} = B_{ji}$, the double-dot product becomes

$$\mathbf{A} = \mathbf{Q} : \mathbf{B} = \left(a\,g_{ij}B_l^{\;l} + d\,B_{ij}\right)\mathbf{b}^i\mathbf{b}^j \; ,$$

depends only on two constants $a$, $d$, having defined $d = b + c$. This is a common condition in continuum mechanics, where some constitutive law links two symmetric rank-2 tensors, like the stress and strain tensors for linear elastic media, or the stress and strain velocity tensors for Newtonian fluids.

# TIME DERIVATIVE OF INTEGRALS OVER MOVING DOMAINS

Some results about time derivatives of integrals over moving domains are collected here. These results are useful for writing balance equations of physical quantities in integral form over arbitrary domains, like:

- integral form of balance equations in continuum mechanics, see Continuum Mechanics:Governing Equations:Integral balance equations for arbitrary domains

- integral form of Maxwell's equations governing classical electromagnetism, see Electromagnetism:Principles of Classical Electromagnetism:derivation of balance equations for arbitrary volume, starting from equations for a control volume

Link to hand-written notes.

## 13.1 Volume density

**Reynolds transport theorem.** Given a volume $V(t)$ with boundary $\partial V(t)$, whose points $\vec{r} \in \partial V(t)$ have velocity $\vec{v}_b$,

$$\frac{d}{dt} \int_{V(t)} f = \int_{V(t)} \frac{\partial f}{\partial t} + \oint_{\partial V(t)} f \vec{v}_b \cdot \hat{n} \, .$$

**"Proof"**

$$\frac{d}{dt} \int_{v(t)} f(\vec{r}, t) \, dt = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \left[ \int_{v(t+\Delta t)} f(\vec{r}, t + \Delta t) - \int_{v(t)} f(\vec{r}, t) \right] =$$

$$= \lim_{\Delta t \to 0} \frac{1}{\Delta t} \left[ \int_{v(t)} \{ f(\vec{r}, t + \Delta t) - f(\vec{r}, t) \} + \int_{\Delta v(t; \Delta t)} f(\vec{r}, t) \right] =$$

$$= \lim_{\Delta t \to 0} \frac{1}{\Delta t} \left[ \int_{v(t)} \left\{ \Delta t \, \frac{\partial f}{\partial t}(\vec{r}, t) + o(\Delta t) \right\} + \oint_{\partial v(t)} \{ \Delta t \, \vec{v}_b \cdot \hat{n} \, f(\vec{r}, t) + o(\Delta t) \} \right] =$$

$$= \int_{v(t)} \frac{\partial f}{\partial t} + \oint_{\partial v(t)} f \vec{v}_b \cdot \hat{n} \, .$$

## 13.2 Flux across a surface

$$\frac{d}{dt}\int_{S(t)}\vec{f}\cdot\hat{n}=\int_{S(t)}\frac{\partial\vec{f}}{\partial t}\cdot\hat{n}+\int_{S(t)}\nabla\cdot\vec{f}\,\vec{v}_b\cdot\hat{n}-\int_{\partial S(t)}\vec{v}_b\times\vec{f}\cdot\hat{t}$$

**"Proof"**

$$\frac{d}{dt}\int_{s(t)}\vec{f}(\vec{r},t)\cdot\hat{n}\,dt=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{s(t+\Delta t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{n}(\vec{r},t+\Delta t)-\int_{s(t)}f(\vec{r},t)\cdot\hat{n}(\vec{r},t)\right]=$$

$$=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{s(t+\Delta t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{n}(\vec{r},t+\Delta t)-\int_{s(t)}f(\vec{r},t)\cdot\hat{n}(\vec{r},t)+\right.$$

$$\left.-\int_{s(t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{n}(\vec{r},t)+\int_{s(t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{n}(\vec{r},t)\right]=$$

$$=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{s(t)}\Delta t\frac{\partial\vec{f}}{\partial t}\cdot\hat{n}+o(\Delta t)+\oint_{\partial\Delta v(t;\Delta t)}\vec{f}\cdot\hat{n}-\oint_{\partial s(t)}\Delta t\vec{f}\cdot\hat{t}\times\vec{v}_b\right]=$$

$$=\int_{s(t)}\frac{\partial\vec{f}}{\partial t}\cdot\hat{n}+\int_s\nabla\cdot\vec{f}\,\vec{v}_b\cdot\hat{n}-\oint_{\partial s(t)}\vec{v}_b\times\vec{f}\cdot\hat{t}\,.$$

having used

$$\oint_{\partial\Delta v(t;\Delta t)}\vec{f}\cdot\hat{n}=\int_{\Delta v}\nabla\cdot\vec{f}=\Delta t\int_s\nabla\cdot\vec{f}\,\hat{n}\cdot\vec{v}_b+o(\Delta t)$$

## 13.3 Work line integral along a line

$$\frac{d}{dt}\int_{\ell(t)}\vec{f}\cdot\hat{t}=\int_{\ell(t)}\frac{\partial\vec{f}}{\partial t}\cdot\hat{t}+\int_{\ell(t)}\nabla\times\vec{f}\cdot\vec{v}_b\times\hat{t}+\vec{f}_B\cdot\vec{v}_B-\vec{f}_A\cdot\vec{v}_A$$

**"Proof"**

$$\frac{d}{dt}\int_{\ell(t)}\vec{f}(\vec{r},t)\cdot\hat{t}(\vec{r},t)\,dt=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{\ell(t+\Delta t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{t}(\vec{r},t+\Delta t)-\int_{\ell(t)}\vec{f}(\vec{r},t)\cdot\hat{t}(\vec{r},t)\right]=$$

$$=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{\ell(t+\Delta t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{t}(\vec{r},t+\Delta t)-\int_{\ell(t)}f(\vec{r},t)\cdot\hat{t}(\vec{r},t)+\right.$$

$$\left.-\int_{\ell(t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{t}(\vec{r},t)+\int_{\ell(t)}\vec{f}(\vec{r},t+\Delta t)\cdot\hat{t}(\vec{r},t)\right]=$$

$$=\lim_{\Delta t\to 0}\frac{1}{\Delta t}\left[\int_{\ell(t+\Delta t)}\Delta t\frac{\partial\vec{f}}{\partial t}(\vec{r},t+\Delta t)\cdot\hat{t}(\vec{r},t+\Delta t)+\right.$$

$$\left.+\oint_{\partial\Delta s(t)}f(\vec{r},t)\cdot\hat{t}(\vec{r},t)-\Delta t\vec{f}_A\cdot\vec{v}_A+\Delta t\vec{f}_B\cdot\vec{v}_B+o(\Delta t)\right]=$$

$$=\int_{\ell(t)}\frac{\partial\vec{f}}{\partial t}\cdot\hat{t}+\int_{\ell(t)}\nabla\times\vec{f}\cdot\vec{v}_b\times\hat{t}+\vec{f}_B\cdot\vec{v}_B-\vec{f}_A\cdot\vec{v}_A\,.$$

having used

$$\oint_{\partial \Delta s(t;\Delta t)} \vec{f} \cdot \hat{t} = \int_{\Delta s} \hat{n} \cdot \nabla \times \vec{f} = \Delta t \int_{\ell} \nabla \times \vec{f} \cdot \hat{v}_b \times \vec{t} + o(\Delta t)$$

**Chapter 13. Time derivative of integrals over moving domains**

# CALCULUS IDENTITIES

Here some calculus identities are listed and proved, often using Cartesian coordinates and index notation. While these identities are independent on the choice of the coordinates - as calculus is - a generic set of coordinates can be used to prove them, and Cartesian coordinates are the most convenint choice.

$$\Delta \vec{v} = \nabla(\nabla \cdot \vec{v}) - \nabla \times \nabla \times \vec{v} \,. \tag{14.1}$$

**Proof.**

Uusing the identity

$$\varepsilon_{ijk}\varepsilon_{ilm} = \delta_{jl}\delta_{km} - \delta_{jm}\delta_{kl}$$

and index notation for Cartesian coordinates,

$$
\begin{aligned}
\left(\nabla \times \nabla \times \vec{v}\right)_i &= \hat{e}_i \cdot \left(\nabla \times \nabla \times \vec{v}\right) = \\
&= \varepsilon_{ijk}\,\partial_j(\varepsilon_{klm}\,\partial_l v_m) = \\
&= \varepsilon_{kij}\,\varepsilon_{klm}\partial_{jl}v_m = \\
&= \left(\delta_{il}\,\delta_{jm} - \delta_{im}\,\delta_{jl}\right)\partial_{jl}v_m = \\
&= \partial_{ij}\,v_j - \partial_{jj}\,v_i = \\
&= \hat{e}_i \cdot \left(\nabla(\nabla \cdot \vec{v}) - \Delta\vec{v}\right) = \qquad = \left(\nabla(\nabla \cdot \vec{v}) - \Delta\vec{v}\right)_i \,.
\end{aligned}
$$

# Part V

# Functional Analysis

# INTRODUCTION TO FUNCTIONAL ANALYSIS

- Lebesgue integral
- $L^p$, $H^p$ function spaces
- Banach and Hilbert spaces

# DIRAC'S DELTA

Dirac's delta $\delta(x)$ is a distribution, or generalized function, with the following properties

1.

$$\int_D \delta(x - x_0)\, dx = 1 \quad \text{if } x_0 \in D \tag{16.1}$$

2.

$$\int_D f(x)\delta(x - x_0)\, dx \quad \text{if } x_0 \in D \tag{16.2}$$

for $\forall f(x)$ "regular" **todo** *what does regular mean?*

## 16.1 Dirac's delta in terms of regular functions

### 16.1.1 Piece-wise constant

$$\delta(x) \sim r_\varepsilon(x) = \begin{cases} \frac{1}{\varepsilon} & x \in \left[-\frac{\varepsilon}{2}, \frac{\varepsilon}{2}\right] \\ 0 & \text{otherwise} \end{cases}$$

**Properties - proof.**

1. Unitariety

$$\int_{x=-\infty}^\infty r_\varepsilon(x - x_0)\, dx = \int_{x=x_0-\frac{\varepsilon}{2}}^{x_0+\frac{\varepsilon}{2}} \frac{1}{\varepsilon}\, dx = 1\,,$$

   for $\forall \varepsilon$;

2. Shift property, using mean-value theorem of continuous functions

$$\int_{x=-\infty}^\infty r_\varepsilon(x - x_0)f(x)\, dx = \int_{x=x_0-\frac{\varepsilon}{2}}^{x_0+\frac{\varepsilon}{2}} \frac{1}{\varepsilon}f(x)\, dx = \frac{1}{\varepsilon}\varepsilon f(\xi)\,,$$

   with $\xi \in \left[x_0 - \frac{\varepsilon}{2}, x_0 + \frac{\varepsilon}{2}\right]$, for the mean value theorem. As $\varepsilon \to 0$, $\xi \to x_0$, and thus

$$\int_{x=-\infty}^\infty r_\varepsilon(x - x_0)f(x)\, dx \to f(x_0)$$

## 16.1.2 Piecewise-linear

$$\delta(x) \sim t_\varepsilon(x) = \begin{cases} \frac{2}{\varepsilon}\left(1 - \frac{2|x|}{\varepsilon}\right) & x \in \left[-\frac{\varepsilon}{2}, \frac{\varepsilon}{2}\right] \\ 0 & \text{otherwise} \end{cases}$$

**Properties - proof**

1. Unitariety

$$\int_{x=-\infty}^{\infty} t_\varepsilon(x - x_0)\, dx = \int_{x=x_0-\frac{\varepsilon}{2}}^{x_0+\frac{\varepsilon}{2}} \frac{2}{\varepsilon}\left(1 - \frac{2|x|}{\varepsilon}\right) dx = \frac{1}{2}\varepsilon\frac{2}{\varepsilon} = 1 \ ,$$

for $\forall \varepsilon$;

2. Shift property, using mean-value integration scheme in $x \in \left[x_0 - \frac{\varepsilon}{2}, x_0\right]$, $x \in \left[x_0, x_0 + \frac{\varepsilon}{2}\right]$ (**todo** *why?*)

$$\int_{x=-\infty}^{\infty} t_\varepsilon(x - x_0)f(x)\, dx = \int_{x=x_0-\frac{\varepsilon}{2}}^{x_0+\frac{\varepsilon}{2}} \frac{2}{\varepsilon}\left(1 - \frac{2|x-x_0|}{\varepsilon}\right) f(x)\, dx =$$

$$= \int_{x=x_0-\frac{\varepsilon}{2}}^{x_0} \frac{2}{\varepsilon}\left(1 - \frac{2|x-x_0|}{\varepsilon}\right) f(x)\, dx + \int_{x=x_0}^{x_0+\frac{\varepsilon}{2}} \frac{2}{\varepsilon}\left(1 - \frac{2|x-x_0|}{\varepsilon}\right) f(x)\, dx =$$

$$= \frac{\varepsilon}{2}\frac{2}{\varepsilon}\left(1 - \frac{2}{\varepsilon}\frac{\varepsilon}{4}\right) f\left(x_0 - \frac{\varepsilon}{4}\right) dx + \frac{\varepsilon}{2}\frac{2}{\varepsilon}\left(1 - \frac{2}{\varepsilon}\frac{\varepsilon}{4}\right) f\left(x_0 + \frac{\varepsilon}{4}\right) dx =$$

$$= \frac{1}{2}f\left(x_0 - \frac{\varepsilon}{4}\right) + \frac{1}{2}f\left(x_0 + \frac{\varepsilon}{4}\right)$$

As $\varepsilon \to 0$

$$\int_{x=-\infty}^{\infty} t_\varepsilon(x - x_0)f(x)\, dx \to f(x_0)$$

## 16.1.3 Gaussian approximation

For $\alpha \to +\infty$,

$$\varphi_\alpha(x) = \sqrt{\frac{\alpha}{\pi}} e^{-\alpha x^2} \sim \delta(x)$$

**Properties - proof**

*Fourier transform* of $\varphi_\alpha(x)$ reads

$$\mathcal{F}\{\varphi_\alpha(x)\}(k) = \int_{x=-\infty}^{+\infty} \varphi_\alpha(x) e^{-ikx}\, dx =$$

$$= \int_{x=-\infty}^{+\infty} \sqrt{\frac{\alpha}{\pi}} e^{-\alpha x^2} e^{-ikx}\, dx =$$

$$= \sqrt{\frac{\alpha}{\pi}} \int_{x=-\infty}^{+\infty} e^{-\alpha\left(x+i\frac{k}{2\alpha}\right)^2}\, dx\, e^{-\frac{k^2}{4\alpha}} =$$

$$= \sqrt{\frac{\alpha}{\pi}} \sqrt{\frac{\pi}{\alpha}}\, e^{-\frac{k^2}{4\alpha}} = e^{-\frac{k^2}{4\alpha}} \ ,$$

for $\alpha \to +\infty$,

$$\mathcal{F}\{\varphi_\alpha(x)\}(k) \to 1$$

Fourier transform of Dirac's delta is 1, as shown in (19.3), thus $\varphi_\alpha(x) \to \delta(x)$ for $\alpha \to +\infty$.

### 16.1.4 Fourier anti-transform

For $a \to +\infty$,

$$\delta(x) \sim \int_{y=-a}^{+a} e^{i2\pi yx} \, dy = \frac{1}{2\pi} \int_{k=-2\pi a}^{2\pi a} e^{ikx} \, dk \ , \tag{16.3}$$

or

$$\delta \sim 2 \int_{y=0}^{a} \cos(2\pi yx) \, dy \ .$$

#### Proof of the equilvanece

$$\delta(x) \sim \frac{1}{2\pi} \int_{k=-2\pi a}^{2\pi a} e^{ikx} \, dk = \frac{1}{2\pi} \left( \int_{k=-2\pi a}^{0} e^{ikx} \, dk + \int_{0}^{k=2\pi a} e^{ikx} \, dk \right) = \frac{1}{2\pi} \int_{k=0}^{2\pi a} \left( e^{ikx} + e^{ikx} \right) \, dx = \frac{1}{\pi} \int_{x=0}^{2\pi a} \cos(kx) \, dk$$

$$= \int_{y=-a}^{+a} e^{i2\pi yx} \, dy = \cdots = \int_{y=0}^{a} \left( e^{i2\pi yx} + e^{i2\pi yx} \right) \, dy = 2 \int_{y=0}^{a} \cos(2\pi yx) \, dy \ .$$

### 16.1.5 sinc$(x)$ approximation

For $a \to +\infty$

$$\delta(x) \sim \frac{\sin(2\pi xa)}{\pi x}$$

#### Proof

Directly follows from integral of the approximation (16.3)

$$\int_{y=-a}^{+a} e^{i2\pi yx} \, dy = \frac{1}{i2\pi x} e^{i2\pi yx} \Big|_{y=-a}^{+a} = \frac{1}{\pi x} \frac{e^{i2\pi ax} - e^{-i2\pi ax}}{2i} = \frac{\sin(2\pi xa)}{\pi x}$$

### 16.1.6 Fourier series

For $x \in [-\pi, \pi]$, and $N \to +\infty$, *Fourier series* of Dirac's delta (train with period $2\pi$) reads

$$\delta(x) \sim \frac{1}{2\pi} \sum_{n=-N}^{N} e^{inx} = \frac{1}{2\pi} \frac{\sin\left(\left(N + \frac{1}{2}\right)x\right)}{\sin\left(\frac{x}{2}\right)}$$

or the $T$-periodic Dirac's delta train,

$$\delta(x) \sim \frac{1}{T} \sum_{n=-N}^{N} e^{in\frac{2\pi}{T}x} \ .$$

**todo** *Write the proof of the last expression, using the relation between complex exponentials and cosine and sine*

### Proof

Coefficients of the Fourier series of Dirac's delta (train with period $T = 2\pi$) are evaluated using the expression (19.2)

$$c_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} \delta(t) e^{-in\frac{2\pi}{2\pi}t} = \frac{1}{2\pi} \ ,$$

and thus the complex Fourier series (19.1) of Dirac's delta reads

$$\delta(x) \sim \sum_{n=-\infty}^{+\infty} c_n e^{in\frac{2\pi}{T}x} = \frac{1}{2\pi} \sum_{n=-\infty}^{+\infty} e^{inx}$$

**Obs.** here, integration interval $[-\pi, \pi]$ to "avoid troubles" with Dirac's delta on the extreme points of the interval (it would give $1/2$ and $1/2$ contributions on both extremes…)

It's possible to write the $T$-periodic Dirac's delta train as

$$\delta(x) \sim \sum_{n=-\infty}^{+\infty} c_n e^{in\frac{2\pi}{T}x} = \frac{1}{T} \sum_{n=-\infty}^{+\infty} e^{in\frac{2\pi}{T}x}$$

**Integral** $I = \int_{-\infty}^{+\infty} e^{-\alpha x^2} \, dx$

$$I^2 = \int_{x=-\infty}^{+\infty} e^{-\alpha x^2} \, dx \int_{y=-\infty}^{+\infty} e^{-\alpha y^2} \, dy =$$

$$= \int_{x=-\infty}^{+\infty} \int_{y=-\infty}^{+\infty} e^{-\alpha(x^2+y^2)} \, dx \, dy =$$

$$= \int_{\theta=0}^{2\pi} \int_{r=0}^{+\infty} e^{-\alpha r^2} r \, dr \, d\theta =$$

$$= 2\pi \frac{1}{2\alpha} \int_{r=0}^{+\infty} e^{-\alpha r^2} \, d\left(\alpha r^2\right) =$$

$$= \frac{\pi}{\alpha} \left[-e^{\alpha r^2}\right]\Big|_{r=0}^{+\infty} = \frac{\pi}{\alpha} \ .$$

# Part VI

# Complex Calculus

# SEVENTEEN

# COMPLEX ANALYSIS

## 17.1 Complex functions, $f : \mathbb{C} \to \mathbb{C}$

A complex function $f$ of complex variable $z = x + iy$, $f : \mathbb{C} \to \mathbb{C}$, can be written as

$$f(z) = \tilde{u}(z) + i\tilde{v}(z) = u(x,y) + iv(x,y) ,$$

as the sum of its real part $u(z)$ and $i$ times its imaginary part $v(x,y)$. Here $x, y \in \mathbb{R}$, while $\tilde{u}(z), \tilde{v}(z) : \mathbb{C} \to \mathbb{R}$ and $u(x,y), v(x,y) : \mathbb{R}^2 \to \mathbb{R}$. With some abuse of notation, tilde won't be always explicitly written when arguments of real and imaginary parts of $f$ functions won't be written.

### 17.1.1 Limit

$$\lim_{z \to z_0} f(z) = f(z_0) \qquad , \qquad \forall \varepsilon > 0 \ \exists \delta > 0 \ \text{ s.t. } \ |f(z) - f(z_0)| < \delta \ \forall z \text{ s.t. } |z - z_0| < \varepsilon, \ z \neq z_0 .$$

### 17.1.2 Derivative

Using the definition of *limit of complex functions*, the derivative of a function $f : \mathbb{C} \to \mathbb{C}$, if it exists, is the limit of incremental ratio,

$$f'(z) = \lim_{\Delta z \to 0} \frac{f(z + \Delta z) - f(z)}{\Delta z} .$$

### 17.1.3 Line Integrals

Given a line $\gamma \in \mathbb{C}$, whose parametric form is $z(s)$, with regular parametrization with parameter $s \in [s_0, s_1]$,

$$\int_\gamma f(z)\, dz = \int_{s=s_0}^{s_1} f(z(s))\, z'(s)\, ds .$$

# 17.2 Holomorphic Functions - Analytic Functions

**Definition 17.2.1**

A holomorphic function is a function whose *derivative* exists.

**Examples of analytic functions. todo**…

## 17.2.1 Cauchy-Riemann conditions

For a holomorphic function $f(z) = u(x,y) + iv(x,y)$, Cauchy-Riemann conditions

$$\begin{cases} u_{/x} = v_{/y} \\ u_{/y} = -v_{/x} \end{cases}$$

hold. The evaluation of the derivative once with $\Delta z = \Delta x$ and once with $\Delta z = i\Delta y$

$$f'(z) = \lim_{\Delta z \to 0} \frac{f(z + \Delta z) - f(z)}{\Delta z} =$$

$$= \begin{cases} \lim_{\Delta x \to 0} \dfrac{f(x + \Delta x, y) - f(x,y)}{\Delta x} = \lim_{\Delta x \to 0} \dfrac{u(x + \Delta x, y) + iv(x + \Delta x, y) - u(x,y) - iv(x,y)}{\Delta x} = u_{/x} + iv_{/x} \\ \lim_{\Delta y \to 0} \dfrac{f(x, y + \Delta y) - f(x,y)}{i\Delta y} = \lim_{\Delta y \to 0} \dfrac{u(x, y + \Delta y) + iv(x, y + \Delta y) - u(x,y) - iv(x,y)}{i\Delta y} = -iu_{/y} + v_{/y} \end{cases}$$

provides the proof.

## 17.2.2 Cauchy Theorem

For a holomorphic function $f$, $f : \Omega \subseteq \mathbb{C} \to \mathbb{C}$

$$\oint_\gamma f(z)\,dz = 0\,,$$

for $\forall \gamma \subset \Omega$. Proof follows from *Green's lemma*, and *Cauchy-Riemann conditions*

$$\oint_\gamma f(z)dz = \oint_\gamma \left(u(x,y) + iv(x,y)\right)(dx + idy) =$$

$$= \oint_\gamma (udx - vdy) + i\oint_\gamma (udy + vdx) =$$

$$= -\int_S \left(\underbrace{u_{/y} + v_{/x}}_{=0}\right) dx\,dy + i\int_S \left(\underbrace{u_{/x} - v_{/y}}_{=0}\right) dx\,dy = 0\,.$$

## 17.3 Useful integrals

### 17.3.1 Independence of line integral for holomorphic functions

For a function $f(z)$ analytic in $D$, the line integral on paths $\ell_{ab,i}$ with the same extreme points $a$, $b$ contained in $D$ is independent on the path, but only depends on the extreme points $a$, $b$,

$$\int_{\ell_{ab,1}} f(z)\, dz = \int_{\ell_{ab,2}} f(z)\, dz$$

The proof readily follows, using *Cauchy theorem* applied to a function $f(z) : D \subseteq \mathbb{C} \to \mathbb{C}$, analytic in $D$, and splitting the closed path $\gamma$ into two paths $\ell_1$, $\ell_2$ with the same extreme points, $\gamma = \ell_1 \cup (-\ell_2)$

$$0 = \oint_\gamma f(z)\, dz = \int_{\ell_1} f(z)\, dz + \int_{-\ell_2} f(z)\, dz = \int_{\ell_1} f(z)\, dz - \int_{\ell_2} f(z)\, dz \, .$$

### 17.3.2 Sum and difference of line integrals

### 17.3.3 Integral of $z^n$

Given a path $\gamma$ embracing $z = 0$ only once in counter-clockwise direction, and $n \in \mathbb{Z}$

$$\oint_\gamma z^n\, dz = \begin{cases} 2\pi i & \text{if } n = -1 \\ 0 & \text{otherwise} \end{cases}$$

Since $z^n$ is analytic everywhere (**todo** *prove it! Add a section with proofs for common functions*) except for $z = 0$, it's possible to evaluate the integral on a circle with center $z = 0$ and radius $R$. Using polar expression of the complex numbers on the circle, $z = Re^{i\theta}$, $\theta \in [0, 2\pi]$, $R$ const, the differential becomes $dz = iRe^{i\theta}d\theta$ and the integral

$$\oint_\gamma z^n\, dz = \int_{\theta=0}^{2\pi} \left(Re^{i\theta}\right)^n iRe^{i\theta}d\theta =$$

$$= i \int_{\theta=0}^{2\pi} R^{n+1} e^{i(n+1)\theta} d\theta =$$

$$= \begin{cases} \text{if } n = -1 & : \quad i2\pi \\ \text{otherwise} & : \quad iR^{n+1} \dfrac{1}{i(n+1)} \left. e^{i(n+1)\theta}\right|_{\theta=0}^{2\pi} = \dfrac{R^{n+1}}{n+1}(1-1) = 0 \end{cases}$$

## 17.4 Meromorphic functions

**Definition 17.4.1**

A meromorphic function in a domain is a function holomorphic everywhere except for a (finite?) number of poles. **check**

## 17.4.1 Singularities

---

**Definition 17.4.2 (Pole)**

A pole of order $n$ of a function $f(z)$ is a complex number $a$ so that

$$f(z) = \frac{\phi(z)}{(z-a)^n} \ ,$$

with $\phi(z)$ holomorphic in $\phi(a) \neq 0$

---

**Examples.** …

---

**Definition 17.4.3 (Branch)**

---

**Examples.** $f(z) = z^{\frac{1}{2}}$

---

**Definition 17.4.4 (Removable singularities)**

---

**Example.** $f(z) = \frac{\sin z}{z}$

**Other irregularities.**

## 17.4.2 Laurent Series

Given a function $f(z)$, in a disk $D_{a,\varepsilon} : 0 < |z-a| < \varepsilon$, its Laurent series centered in $a$ is the convergent (to $f(z)$, **todo** *which type of convergnence?*) series

$$f(z) \sim \sum_{n=-\infty}^{+\infty} a_n (z-a)^n \ , \tag{17.1}$$

with

$$a_n = \frac{1}{2\pi i} \int_\gamma f(z)\,(z-a)^{-(n+1)}\,dz \tag{17.2}$$

and $\gamma$ embracing $z = a$ once counter-clockwise. Proof follows immediately inserting the expressions of the coefficients $a_n$ and using the *integral of $z^n$*. Evaluating the integral (17.2) of the coefficients of the Laurent series, using (17.1) to replace $f(z)$ with its series

$$a_n = \frac{1}{2\pi i} \oint_\gamma \sum_{m=-\infty}^{+\infty} a_m (z-a)^m (z-a)^{-(n+1)} =$$

$$= \frac{1}{2\pi i} \oint_\gamma \sum_{m=-\infty}^{+\infty} a_m (z-a)^{m-n-1}\,dz =$$

$$= \frac{1}{2\pi i} \oint_\gamma a_n\,z^{-1}\,dz =$$

$$= a_n \ .$$

**todo** *Some freestyle with function and its convergent series…add some detail, and the meaning of convergence*

---

### 17.4.3 Cauchy formula

For an analytic function $f(z)$,

$$f(a) = \frac{1}{2\pi i} \oint_\gamma \frac{f(z)}{z-a} \, dz$$

Proof readily follows using the *integral of $z^n$* on the Taylor series of $\frac{f(z)}{z-a}$ whose $0^{th}$ order term reads $f(a)$,

$$\frac{1}{2\pi i} \oint_\gamma \frac{f(a) + \sum_{m=1}^{+\infty} f'(a)(z-a)^m}{z-a} \, dz = \frac{1}{2\pi i} \oint_\gamma \frac{f(a)}{z-a} \, dz = f(a)\frac{2\pi i}{2\pi i} = f(a) \,.$$

### 17.4.4 Residues

---

**Definition 17.4.5 (Residue)**

The residue of function $f$ in $a$, $\mathrm{Res}(f,a)$ is a complex number $R$ so that $f(z) - \frac{R}{(z-a)}$ has analytic antiderivative in a disk $D_{a,\varepsilon} : 0 < |z-a| < \varepsilon$.

---

**todo** Explain this definition. Couldn't be possible to use $\mathrm{Res}(f,a) = \frac{1}{2\pi i} \oint_\gamma f(z) \, dz = a_{-1}$ instead?

**Properties.**

- If $f(z)$ is analytic in $D_{a,\varepsilon}$ and has a pole of order $n$ in $z = a$, its Laurent series has $a_m = 0$ for $m < n$ and reads

$$f(z) = \sum_{m=-n}^{+\infty} a_m (z-a)^m \,, \tag{17.3}$$

with $a_{-n} \neq 0$. Since $f(z)$ has a pole of order $n$ in $z = a$, it can be written as

$$f(z) = \frac{\phi(z)}{(z-a)^n} \,,$$

with $\phi(z)$ analytic in $D_{a,\varepsilon}$ and $\phi(a) \neq 0$. Since $\phi(z)$ is analytic, it has a Taylor series (or a Laurent series with non-negative powers),

$$\phi(z) \sim \sum_{m=0}^{+\infty} b_m (z-a)^m \,,$$

(**todo** *prove it! Extension of the real case. Add a link to the proof*) and thus

$$f(z) \sim \sum_{m=0}^{+\infty} b_m (z-a)^{m-n} = \sum_{m=-n}^{+\infty} b_{m+n}(z-a)^m = \sum_{m=-n}^{+\infty} a_m(z-a)^m \,,$$

with $a_m = b_{m+n}$.

- For simple closed path $\gamma$ (embracing $a$ only once counter-clokwise) in $D_{a,\varepsilon}$,

$$\oint_\gamma f(z) \, dz = 2\pi i a_{-1} = 2\pi i \mathrm{Res}(f,a) \tag{17.4}$$

The proof readily follows, using the *integral of $z^n$* and Laurent series (17.1) of $f(z)$,

$$\oint_\gamma f(z) \, dz = \oint_\gamma \sum_{m=-\infty}^{+\infty} a_m (z-a)^m \, dz = 2\pi i a_{-1} \,.$$

- For a pole $a$ of order $n$, the following holds

$$a_{-1} = \frac{1}{(n+1)!} \lim_{z \to a} \frac{d^{n-1}}{dz^{n-1}} \left[ (z-a)^n f(z) \right]$$

The proof follows using Laurent series {eq}`eq:laurent:pole-n` for a function with pole of order $n$, and evaluating the $(n-1)^{th}$ order derivative

$$
\frac{d^{n-1}}{dz^{n-1}} \left[ (z-a)^n f(z) \right] = \frac{d^{n-1}}{dz^{n-1}} \left[ (z-a)^n \sum_{m=-n}^{+\infty} a_n (z-a)^m \right] =
$$

$$
= \frac{d^{n-1}}{dz^{n-1}} \left[ \sum_{m=-n}^{+\infty} a_n (z-a)^{m+n} \right] =
$$

$$
= \frac{d^{n-1}}{dz^{n-1}} \left[ \sum_{m=0}^{+\infty} a_{m-n} (z-a)^m \right] =
$$

$$
= \frac{d^{n-2}}{dz^{n-2}} \left[ \sum_{m=0}^{+\infty} m\, a_{m-n} (z-a)^{m-1} \right] =
$$

$$
= \frac{d^{n-3}}{dz^{n-3}} \left[ \sum_{m=0}^{+\infty} m(m-1) a_{m-n} (z-a)^{m-2} \right] =
$$

$$
= \cdots =
$$

$$
= \left[ \sum_{m=0}^{+\infty} m!\, a_{m-n} (z-a)^{m-n+1} \right]
$$

and then letting $z \to a$, so that only the term with $m - n + 1 = 0$ survives

$$
\lim_{z \to a} \frac{d^{n-1}}{dz^{n-1}} \left[ (z-a)^n \sum_{m=-n}^{+\infty} a_n (z-a)^m \right] = (n-1)!\, a_{-1} \ .
$$

### 17.4.5 Residue Theorem

**Theorem 17.4.1 (Residue Theorem)**

Given $f(z)$ with a finite number of poles $p_n \in D$, then

$$
\int_\gamma f(z)\, dz = 2\pi i \sum_n I(\gamma, p_n) \text{Res}(f, p_n) \ ,
$$

being $\gamma$ a path in $D$, and $I(\gamma, p_n)$ the winding index of the path $\gamma$ around pole $p_n$ (+1 for each counter-clockwise loop, -1 for each clockwise loop).

The proof readily follows extending the result for a single pole (17.4) to general number of poles and general paths $\gamma$ embracing (with sign) each pole $p_n$ $I(\gamma, p_n)$ times, with the same techinques shown in section *Sum and difference of line integrals*.

### 17.4.6 Cauchy argument principle

### 17.4.7 Evaluation of integrals

### 17.4.8 Inverse Laplace Transform

Given Laplace transform

$$F(s) := \mathcal{L}\{f(t)\}(s) := \int_{t=0^-}^{+\infty} f(t)e^{-st}\,dt\,,$$

the inverse transform can be evaluated as

$$f(t) = \mathcal{L}^{-1}\{F(s)\}(t) := \lim_{T\to+\infty} \frac{1}{2\pi i}\int_{s=a-iT}^{a+iT} e^{st}F(s)\,ds\,,$$

with $a > \mathrm{Re}\{p_n\}$ (**todo** *why?*) for each pole of the function $F(s)$, evaluated on the vertical line $s = a + iy$, $y \in [-T, T]$, $ds = idy$,

$$
\begin{aligned}
\lim_{T\to+\infty} \frac{1}{2\pi i}\int_{s=a-iT}^{a+iT} e^{st}F(s)\,ds &= \lim_{T\to+\infty} \frac{1}{2\pi i}\int_{s=a-iT}^{a+iT} e^{st}\int_{\tau=0^-}^{+\infty} f(\tau)e^{-s\tau}\,d\tau\,ds = \\
&= \lim_{T\to+\infty} \frac{1}{2\pi i}\int_{y=-T}^{T} e^{(a+iy)t}\int_{\tau=0^-}^{+\infty} f(\tau)e^{-(a+iy)\tau}\,d\tau\,idy = \\
&= \lim_{T\to+\infty} \frac{1}{2\pi}\int_{y=-T}^{T}\int_{\tau=0^-}^{+\infty} e^{iy(t-\tau)}e^{a(t-\tau)}f(\tau)\,d\tau\,dy = \\
&= \dots \\
&= \int_{\tau=0^-}^{+\infty} \delta(t-\tau)e^{a(t-\tau)}f(\tau)d\tau = f(t)\,.
\end{aligned}
$$

having used the transform of *Dirac's delta* $\delta(t) = \frac{1}{2\pi}\int_{\omega=-\infty}^{+\infty} e^{-j\omega t}\,d\omega$.

**todo** *Ohter approach: if $a > Re\{p_n\}$, the contour built with the vertical line with real part $a$ and the arc of circumference on its...*

# LAPLACE TRANSFORM

*Definition of Laplace transform.* Definition of Laplace transform, inverse transform, properties and theorems.

*Applications of Laplace transform.* Solution of ODEs,…

## 18.1 Definition and Properties

### 18.1.1 Definition of Lapalce transform

$$\mathcal{L}\left\{f(t)\right\}(s) := \int_{t=0^-}^{+\infty} e^{-st} f(t)\, dt = F(s)\,.$$

### 18.1.2 Inverse transform

$$f(t) = \mathcal{L}^{-1}\left\{F(s)\right\} = ...$$

### 18.1.3 Properties

**Linearity.**

$$\mathcal{L}\{af(t) + bg(t)\}(s) = aF(s) + bG(s)$$

*Dirac delta.*

$$\mathcal{L}\left\{\delta(t)\right\} = \int_{t=0^-}^{+\infty} \delta(t)\, e^{st}\, dt = 1$$

**Time delay.** If $f(t) = 0$ for $t < 0$ ("causality"), for $\tau > 0$,

$$\mathcal{L}\{f(t-\tau)\}(s) = e^{-s\tau} F(s)$$

Proof readily follows direct computation with change of variable $z = t - \tau$, $dt = dz$

$$\mathcal{L}\{f(t-\tau)\}(s) = \int_{t=0^-}^{+\infty} f(t-\tau)e^{-st}\, dt = \int_{z=-\tau}^{+\infty} f(z)e^{-sz}\, dz\, e^{-s\tau} = \int_{z=0}^{+\infty} f(z)e^{-sz}\, dz\, e^{-s\tau} = e^{-s\tau} F(s)\,.$$

**"Frequency shift"**

$$\mathcal{L}\{f(t)e^{at}\}(s) = F(s-a)$$

Direct computation gives

$$\mathcal{L}\{f(t)e^{at}\}(s) = \int_{t=0^-}^{+\infty} f(t)e^{at}e^{-st}\,dt = \int_{t=0^-}^{+\infty} f(t)e^{-(s-a)t}\,dt = F(s-a)$$

**Derivative.**

$$\mathcal{L}\{f'(t)\}(s) = sF(s) - f(0^-)\,.$$

Proof readily follows direct computation, with integration by parts

$$\mathcal{L}\{f'(t)\}(s) = \int_{t=0^-}^{+\infty} f'(t)e^{-st}\,dt = [f(t)e^{-st}]\,|_{t=0^-}^{+\infty} + s\int_{t=0^-}^{+\infty} f(t)e^{-st}\,dt = sF(s) - f(0^-)\,,$$

provided that $\lim_{s\to+\infty} f(t)e^{-st} = 0$.

**Integral.**

$$\mathcal{L}\left\{\int_{\tau=0}^{t} f(\tau)\,d\tau\right\}(s) = \frac{1}{s}F(s)\,.$$

Proof readily follows direct computation, with integration by parts

$$\mathcal{L}\left\{\int_{\tau=0^-}^{t} f(\tau)\,d\tau\right\}(s) = \int_{t=0^-}^{+\infty}\int_{\tau=0^-}^{t} f(\tau)\,d\tau e^{-st}\,dt = \left[-\frac{e^{-st}}{s}\int_{\tau=0^-}^{t} f(\tau)\,d\tau\right]_{t=0}^{+\infty} + \frac{1}{s}\int_{t=0}^{+\infty} f(t)e^{-st}\,dt = \frac{1}{s}F(s)\,,$$

provided that $\int_{\tau=0^-}^{0} f(\tau)d\tau = 0$ and $\lim_{t\to+\infty} \frac{e^{-st}}{s}\int_{\tau=0^-}^{+\infty} f(\tau)\,d\tau = 0$.

**Convolution.**

$$\begin{aligned}
\mathcal{L}\{f(t)*g(t)\} &= \int_{t=0^-}^{+\infty}\int_{\tau=-\infty}^{+\infty} f(t-\tau)g(\tau)\,d\tau\,e^{-st}\,dt = && (1) \\
&= \int_{\tau=-\infty}^{+\infty}\int_{z=-\tau^-}^{+\infty} f(z)g(\tau)\,e^{-s(z+\tau)}\,d\tau\,dz = && (2) \\
&= \int_{z=0^-}^{+\infty} f(z)\,e^{-sz}\,dz\int_{\tau=0^-}^{+\infty} g(\tau)e^{-s\tau} = \\
&= \mathcal{L}\{f(t)\}(s)\,\mathcal{L}\{g(t)\}(s)\,.
\end{aligned} \qquad (18.1)$$

having performed the change of coordinates $z = t - \tau, \tau = \tau$, with unitary Jacobian,

$$\frac{\partial(t,\tau)}{\partial(z,\tau)} = \partial_z t\partial_\tau\tau - \partial_z\tau\partial_\tau t = 1\cdot 1 - 1\cdot 0 = 1,$$

given the proper description of the domain of integration summarised in the extremes of integration in (1), and causality - i.e. all the functions $f(t)$ are identically zero for $t < 0$ - in (2).

**Initial value.** If …

$$f(0^+) = \lim_{s\to+\infty} sF(s)$$

From direct computation,

$$\begin{aligned}
\lim_{s\to+\infty} sF(s) &= \lim_{s\to+\infty} s\int_{t=0^-}^{+\infty} f(t)\,e^{-st}\,dt = \\
&= \lim_{s\to+\infty}\left\{\left[s\left(-\frac{e^{-st}}{s}\right)f(t)\right]\Big|_{t=0}^{+\infty} + \int_{t=0}^{+\infty} e^{-st}f'(t)\,dt\right\} = \\
&= \lim_{s\to+\infty}\left\{[-e^{-st}f(t)]\Big|_{t=0}^{+\infty} + \int_{t=0}^{+\infty} e^{-st}f'(t)\,dt\right\} = \\
&= f(0)\,,
\end{aligned}$$

provided that $\lim_{s \to +\infty} \lim_{t \to +\infty} e^{-st} f(t) = 0$ and $\lim_{s \to +\infty} \int_{t=0}^{+\infty} e^{-st} f'(t)\, dt = 0$.

**Final value.** If …

$$f(+\infty) = \lim_{s \to 0} sF(s)$$

From direct computation (**todo** *check and/or explain proof*),

$$\lim_{s \to 0} sF(s) = \lim_{s \to 0} s \int_{t=0^-}^{+\infty} f(t)\, e^{-st}\, dt =$$

$$= \lim_{s \to 0} \left\{ \left[ s \left( -\frac{e^{-st}}{s} \right) f(t) \right] \Big|_{t=0}^{+\infty} + \int_{t=0}^{+\infty} e^{-st} f'(t)\, dt \right\} =$$

$$= \lim_{s \to 0} \left\{ \left[ -e^{-st} f(t) \right] \Big|_{t=0}^{+\infty} + \int_{t=0}^{+\infty} e^{-st} f'(t)\, dt \right\} =$$

$$= f(0) + f(+\infty) - f(0) = f(+\infty)\ ,$$

provided that $\lim_{s \to 0} \lim_{t \to +\infty} e^{-st} f(t) = 0$.

# 18.2  Applications of Laplace Transform

## 18.2.1  Ordinary differential equations

Exploiting the properties of transforming derivation into product by the complex variable $s$, Laplace transform can be a useful tool in solving *Ordinary differential equations*.

Linear ODEs are treated in details in the section about *Linear Time-Invariant Systems*.

### First-order linear differential equation with constant coefficients

A Cauchy problem governed by a first-order linear differential equation with constant coefficients reads

$$\begin{cases} \dot{u} + au = f(t) \\ \text{for } t > 0 \\ u(0^-) = u_0\ . \end{cases}$$

**Comments.**  1) For linear differential equations with constant coefficients, a time shift is always possible to set initial conditions in $t_0 = 0$; 2) dealing with **impulsive forces** (more on this in *LTI systems: impulsive forces*), some attention needs to be paid to the time of initial conditions: if impulsive forces at $t = 0$ may exist, initial conditions need to be set at $t = 0^-$; the effect of the impulsive force is equivalent to a jump in initial conditions from $t = 0^-$ and $t = 0^+$.

### General solution in time domain

Multiplying the ODE by $e^{at}$,

$$e^{at} f(t) = e^{at} \left( \dot{u} + au \right) = \frac{d}{dt} \left( e^{at} u(t) \right)$$

and integrating from $0^-$ to a generic time $t$ — after changing the dummy integration variable from $\tau$ to $t$ —

$$e^{at} u(t) - u(0^-) = \int_{\tau=0^-}^{t} e^{a\tau} f(\tau) d\tau\ ,$$

and thus

$$u(t) = e^{-at}u(0^-) + \int_{\tau=0^-}^{t} e^{-a(t-\tau)}f(\tau)d\tau \;.$$

The solution of the problem can be also computed using Laplace transform:

1. transforming the problem from time to Laplace domain: the differential problem with unknown $u(t)$ in time is transformed into an algebraic problem $\hat{u}(s)$. Using Laplace **transform of the time derivative**

$$s\hat{u}(s) + u_0 + a\hat{u}(s) = \hat{f}(s) \;,$$

2. solving the algebraic problem in Laplace domain for $\hat{u}(s)$

$$\hat{u}(s) = \frac{1}{s+a}u_0 + \frac{1}{s+a}f(s) \;.$$

3. transforming back the solution in time domain, $u(t) = \mathcal{L}^{-1}\{\hat{u}(s)\}$,

$$u(t) = \mathcal{L}^{-1}\{\hat{u}(s)\}(t) =$$
$$= \mathcal{L}^{-1}\left\{\frac{u_0}{s+a}\right\} + \mathcal{L}^{-1}\left\{\frac{1}{s+a}\hat{f}(s)\right\} =$$
$$= e^{-at}u_0 + \int_{0^-}^{t} e^{-a(t-\tau)}f(\tau)d\tau \;,$$

having used the inverse transform of the exponential $\mathcal{L}\{e^{ct}\} = \int_{t=0^-}^{+\infty} e^{ct}e^{-st}\,dt = \frac{1}{s-c}$, and the inverse transform of the convolution

$$\hat{f}(s)\hat{g}(s) = \mathcal{L}\{f(t) * g(t)\}(s) = \int_{t=0^-}^{+\infty}\left\{\int_{\tau=0^-}^{t} f(\tau)g(t-\tau)d\tau\right\}e^{-st}dt \;.$$

# FOURIER TRANSFORMS

Fourier transforms are linear transformations of functions usually relating a physical domain of time and/or space, with a domain of frequency and/or wave-vectors.

Fourier transforms can be useful in:

- highlighting the frequency content of functions

- solving problems: sometimes, it can be easier to transform a problem in frequency domain, solve it in frequency domain, and transform the solution back to the physical domain

**Contents.**

*Fourier series*. Fourier series is defined for finite-domain or periodic, time-continuous functions, or - more generally - continuous functions in the physical domain.

*Fourier transform*. Fourier transform is defined for infinite-domain non-periodic, time-continuous functions, or - more generally - continuous functions in the physical domain.

*Relations between Fourier transforms and sampling*. Fourier series, Fourier transform, discrete time Fourier transform and discrete Fourier transforms are presented, and their relations discussed. Fundamental results about **evenly-spaced sampling** seamlessly follows, as **Shannon-Nyquist theorem**, *Theorem 19.3.1*, shows.

Different Fourier transforms exist, depending if the original function is:

- time discrete/time continuous

- periodic/non-periodic

## 19.1 Fourier Series

For a $T$-periodic function,

$$g(t) \sim \frac{a_0}{2} + \sum_{n=1}^{+\infty} \left[ a_n \, \cos\left( n\frac{2\pi}{T}t \right) + b_n \, \sin\left( n\frac{2\pi}{T}t \right) \right] ,$$

**todo** Prove it with properties of integrals of sin and cos over $t \in [0, T]$; prove convergence to average value at jumps

The exponential form reads

$$g(t) \sim \sum_{n=-\infty}^{+\infty} c_n e^{in\frac{2\pi}{T}t} , \tag{19.1}$$

where

$$c_n = \frac{1}{T} \int_{t=0}^{T} f(t)\, e^{-in\frac{2\pi}{T}t} \, . \tag{19.2}$$

### Proof

Exploiting the properties of integrals of complex exponentials with $k \in \mathbb{Z}$

$$\int_{t=0}^{T} e^{ik\frac{2\pi}{T}t}\, dt = \begin{cases} \frac{1}{ik\frac{2\pi}{T}} \left[ e^{ik\frac{2\pi}{T}t} \right]\Big|_{t=0}^{T} = 0 & \text{if } k \neq 0 \\ T & \text{if } k = 0 \end{cases}$$

$$\int_{t=0}^{T} f(t) e^{-im\frac{2\pi}{T}t}\, dt \sim \int_{t=0}^{T} \sum_{n=-\infty}^{+\infty} c_n e^{in\frac{2\pi}{T}t} e^{-im\frac{2\pi}{T}t} \sim \sum_{n=-\infty}^{+\infty} c_n \sim \int_{t=0}^{T} e^{i(n-m)\frac{2\pi}{T}t} \sim T\, c_m \, .$$

## 19.2 Fourier Transform

**Contents**: *definition*; *properties*; *inverse transform*; *Plancherel's theorem*; *uncertainty relation*

### 19.2.1 Definition

…

$$\mathcal{F}\{g(t)\}(f) := \int_{t=-\infty}^{+\infty} g(t)\, e^{-i2\pi ft} \, dt.$$

### 19.2.2 Properties

**Linearity**

*Dirac delta*.

$$\mathcal{L}\{\delta(t)\} = \int_{t=-\infty}^{+\infty} \delta(t)\, e^{-i2\pi ft}\, dt = 1 \tag{19.3}$$

**Time delay.**

**Derivative.**

**Integral.**

**Initial value.**

**Final value.**

---

**Property 19.2.1 (Transform of convolution)**

$$\mathcal{F}\{a * b(t)\}(f) = \mathcal{F}\{a(t)\}(f)\, \mathcal{F}\{b(t)\}(f)$$

---

**Proof.**

$$\mathcal{F}\left\{a * b(t)\right\}(f) = \int_{t=-\infty}^{+\infty}\int_{\tau=-\infty}^{+\infty} a(\tau)b(t-\tau)\, d\tau\, e^{-i2\pi ft}\, dt = \qquad (1)$$

$$= \int_{\tau=-\infty}^{+\infty} a(\tau)\, e^{-i2\pi f\tau}\, d\tau \int_{z=-\infty}^{+\infty} b(z)\, e^{-i2\pi fz} =$$

$$= \mathcal{F}\left\{a(t)\right\}(f)\,\mathcal{F}\left\{b(t)\right\}(f)\,,$$

having used (1) transformation of coordinates $(z,\tau) = (t-\tau,\tau)$ with unit Jacobian

$$\frac{\partial(z,\tau)}{\partial(t,\tau)} = \begin{vmatrix} 1 & -1 \\ 0 & 1 \end{vmatrix} = 1\,.$$

### 19.2.3 Inverse Fourier Transform

Under the assumptions …**todo**, the inverse Fourier transform reads

$$\mathcal{F}^{-1}\left\{G(f)\right\}(t) := \int_{f=-\infty}^{+\infty} G(f)\, e^{i2\pi ft}\, df.$$

**Proof using Dirac's delta expression.**

$$\mathcal{F}^{-1}\left\{G(f)\right\}(t) := \int_{f=-\infty}^{+\infty} G(f)\, e^{i2\pi ft}\, df = \int_{f=-\infty}^{+\infty}\int_{\tau=-\infty}^{+\infty} g(\tau)e^{-i2\pi f\tau}\, e^{i2\pi ft}\, df =$$

$$= \int_{f=-\infty}^{+\infty}\int_{\tau=-\infty}^{+\infty} g(\tau)e^{-i2\pi f\tau}\, e^{i2\pi ft}\, df =$$

$$= \int_{f=-\infty}^{+\infty}\int_{\tau=-\infty}^{+\infty} g(\tau)e^{i2\pi f(t-\tau)}\, df =$$

$$= \int_{\tau=-\infty}^{+\infty} g(\tau)\delta(t-\tau)\, d\tau = g(t)\,.$$

**Proof using dominated convergence theorem and Fubini's lemma.**

**Proof.** By the *dominated convergence theorem*, it follows that

$$\int_{\mathbb{R}} e^{i2\pi x\xi} F(\xi)\, d\xi = \lim_{\varepsilon\to 0}\int_{\mathbb{R}} \underbrace{e^{-\pi\varepsilon^2\xi^2+i2\pi x\xi}}_{G(\xi;x,\varepsilon)} F(\xi)\, d\xi =$$

$$= \lim_{\varepsilon\to 0}\int_{\mathbb{R}} g(y;x,\varepsilon)f(y)\, dy =$$

$$= \lim_{\varepsilon\to 0}\int_{\mathbb{R}} \varphi_\varepsilon(x-y)\, f(y)\, dy =$$

$$= \int_{\mathbb{R}} \delta(x-y)\, f(y)\, dy = f(x)$$

**Lemma 1.** The Fourier transform of function $\varphi(t) := e^{-\pi|t|^2}$ reads

$$
\mathcal{F}\{\varphi(t)\}(\omega) = \int_{t=-\infty}^{+\infty} \varphi(t)e^{-i\omega t}\,dt =
$$

$$
= \int_{t=-\infty}^{+\infty} e^{-\pi|t|^2}e^{-i\omega t}\,dt =
$$

$$
= \int_{t=-\infty}^{+\infty} e^{-\pi\left(t^2+i\frac{\omega}{\pi}t-\frac{\omega^2}{4\pi^2}\right)}\,dt\, e^{-\frac{\omega^2}{4\pi^2}} =
$$

$$
= \int_{t=-\infty}^{+\infty} e^{-\pi\left(t+i\frac{\omega}{2\pi}\right)^2}\,dt\, e^{-\frac{\omega^2}{4\pi}} =
$$

$$
= e^{-\frac{\omega^2}{4\pi}},
$$

having evaluated *the integral* $\int_{-\infty}^{+\infty} e^{-\alpha x^2}$ with $\alpha = \pi$. **todo** *justify the result for complex exponential. Use Bromwich contour integrals*

**Lemma 2.** Fourier transform of $f(\alpha t),\, \alpha > 0$

$$
\mathcal{F}\{f(\alpha t)\}(\omega) = \int_{\mathbb{R}} f(\alpha t)e^{-j\omega t}\,dt = \int_{\tau\in\mathbb{R}} f(\tau)e^{-j\frac{\omega}{\alpha}\tau}\,d\tau\frac{1}{\alpha} = \frac{1}{\alpha}F\left(\frac{\omega}{\alpha}\right)
$$

**Lemma 3.** $\frac{1}{\varepsilon}\varphi\left(\frac{t}{\varepsilon}\right) \to \delta(x)$ for $\varepsilon \to 0$

$$
\mathcal{F}\left\{\frac{1}{\varepsilon}\varphi\left(\frac{t}{\varepsilon}\right)\right\}(\omega) = \frac{1}{\varepsilon}\varepsilon e^{-\frac{\omega^2}{4\pi\varepsilon^2}} = e^{-\frac{\omega^2}{4\pi\varepsilon^2}}
$$

0. Fourier transform

$$
G(f) = \int_{t=-\infty}^{\infty} e^{-i\omega t}g(t)\,dt
$$

1.

$$
g(t) = e^{i\alpha t}\psi(t)
$$

$$
\mathcal{F}\{g(t)\}(\omega) = \int_{t=-\infty}^{+\infty} g(t)e^{-i\omega t}\,dt = \int_{t=-\infty}^{+\infty} \psi(t)e^{i\alpha t}e^{-i\omega t}\,dt = \int_{t=-\infty}^{+\infty} \psi(t)e^{-i(\omega-\alpha)t}\,dt = \mathcal{F}\{\psi(t)\}(\omega - \alpha)\,.
$$

2.

$$
\psi(t) = \phi(\alpha t)
$$

$$
\mathcal{F}\{\psi(t)\} = \int_{t=-\infty}^{+\infty} \psi(t)e^{-i\omega t}\,dt = \int_{t=-\infty}^{+\infty} \phi(\alpha t)e^{-i\omega t}\,dt = \int_{\tau=-\infty}^{+\infty} \phi(\tau)e^{-i\frac{\omega}{\alpha}\tau}\frac{d\tau}{\alpha} = \frac{1}{\alpha}\mathcal{F}\{\phi(t)\}\left(\frac{\omega}{\alpha}\right)\,.
$$

3. Fubini's theorem

4.

$$
\varphi(t) := e^{-\pi t^2}
$$

$$
\mathcal{F}\{\varphi(t)\} = \int_{t=-\infty}^{+\infty} \varphi(t)e^{-i\omega t}\,dt = \int_{t=-\infty}^{+\infty} e^{-\pi t^2}e^{-i\omega t}\,dt
$$

$$
0 = \oint_{\gamma} e^{-\alpha|z|^2}\,dz = \int_{...} ...
$$

$$
z = Re^{i\theta}, \quad dz = iRe^{i\theta}\,d\theta
$$

$$
\int_{C/4} e^{-\alpha|z|^2}\,dz = \int_{\theta=0}^{\frac{\pi}{2}} e^{-\alpha R^2}iRe^{i\theta}\,d\theta = iRe^{-\alpha R^2}\frac{e^{-i\theta}}{i}\Big|_{\theta=0}^{\frac{\pi}{2}}
$$

$$\int_{t=0}^{+\infty} e^{-\pi t^2} e^{-i\omega t}\,dt = \int_{t=0}^{+\infty} e^{-\left(\pi t^2 + i\omega t - \frac{\omega^2}{4\pi}\right)}\,dt\,e^{-\frac{\omega^2}{4\pi}} =$$

$$= \int_{t=0}^{+\infty} e^{-\pi\left(t + i\frac{\omega}{2\pi}\right)^2}\,dt\,e^{-\frac{\omega^2}{4\pi}}$$

5. $\varphi_\varepsilon(t) = \frac{1}{\varepsilon^n}\varphi\left(\frac{t}{\varepsilon}\right)$, $t \in \mathbb{R}^n$, is an approximation of Dirac's delta for $\varepsilon \to 0$, so that

$$\lim_{\varepsilon \to 0} \int_{t=-\infty}^{+\infty} \varphi_\varepsilon(t-\tau)f(t)\,dt = f(\tau)$$

$$\lim_{\varepsilon \to 0} \int_{t=-\infty}^{+\infty} \varphi_\varepsilon(t)\,dt = 1$$

As the Fourier transform $\mathcal{F}\left\{\varphi_\varepsilon(t)\right\}(\omega) \to 1$ for $\varepsilon \to 0$, then $\varphi_\varepsilon(t) \to \delta(t)$.

## 19.2.4 Plancherel's theorem

…assumptions…**todo**

$$\int_{f=-\infty}^{+\infty} |G(f)|^2\,df = \int_{t=-\infty}^{+\infty} |g(t)|^2\,dt \tag{19.4}$$

and

$$\int_{f=-\infty}^{+\infty} A^*(f)\,G(f)\,df = \int_{t=-\infty}^{+\infty} a^*(t)\,g(t)\,dt\ . \tag{19.5}$$

### Proof of Plancherel's thm for the magnitude

$$\int_{f=-\infty}^{+\infty} |G(f)|^2\,df = \int_{f=-\infty}^{+\infty} G(f)^*G(f)\,df =$$

$$= \int_{f=-\infty}^{+\infty} \left(\int_{t_1=-\infty}^{+\infty} g(t_1)e^{-i2\pi f t_1}dt_1\right)^* \left(\int_{t_2=-\infty}^{+\infty} g(t_2)e^{-i2\pi f t_2}dt_2\right)\,df =$$

$$= \int_{t_1,t_2=-\infty}^{+\infty} g^*(t_1)\,g(t_2) \int_{f=-\infty}^{+\infty} e^{i2\pi f(t_1-t_2)}\,df\,dt_1\,dt_2 = \tag{1}$$

$$= \int_{t_1,t_2=-\infty}^{+\infty} g^*(t_1)\,g(t_2)\,\delta(t_1-t_2)\,dt_1\,dt_2 = \tag{2}$$

$$= \int_{t_1=-\infty}^{+\infty} g^*(t_1)\,g(t_1)\,dt_1 = \tag{3}$$

$$= \int_{t_1=-\infty}^{+\infty} |g(t_1)|^2\,dt_1\ .$$

having used (1) the approximation (16.3) of Dirac's delta, and (2) property (16.2) of Dirac's delta, and (3) the expression of the absolute value of complex functions $g^*(t_1)g(t_1) = |g(t_1)|^2$.

**Proof of Plancherel's thm for the product of functions**

$$\int_{f=-\infty}^{+\infty} A^*(f)\, G(f)\, df = \int_{f=-\infty}^{+\infty} G(f)^* G(f)\, df =$$

$$= \int_{f=-\infty}^{+\infty} \left( \int_{t_1=-\infty}^{+\infty} a(t_1) e^{-i2\pi f t_1}\, dt_1 \right)^* \left( \int_{t_2=-\infty}^{+\infty} g(t_2) e^{-i2\pi f t_2}\, dt_2 \right)\, df =$$

$$= \int_{t_1,t_2=-\infty}^{+\infty} a^*(t_1)\, g(t_2) \int_{f=-\infty}^{+\infty} e^{i2\pi f(t_1-t_2)}\, df\, dt_1\, dt_2 = \qquad (1)$$

$$= \int_{t_1,t_2=-\infty}^{+\infty} a^*(t_1)\, g(t_2)\, \delta(t_1-t_2)\, dt_1\, dt_2 = \qquad (2)$$

$$= \int_{t_1=-\infty}^{+\infty} a^*(t_1)\, g(t_1)\, dt_1 \,.$$

having used (1) the approximation (16.3) of Dirac's delta, and (2) property (16.2) of Dirac's delta.

## 19.2.5 Uncertainty relation

An uncertainty relation holds linking standard deviations of a probability density function in time domain and a probability density function built with its Fourier transform. From this very same relation, Heisenberg uncertainty relation between position and momentum in Quantum Mechanics seamlessly follows.

Given a function $g(t)$ whose square of the absolute value is normalized to one, and thus it can be used as a probability density function in time domain,

$$\int_{t=-\infty}^{+\infty} |g(t)|^2\, dt = \int_{t=-\infty}^{+\infty} g^*(t)\, g(t)\, dt = 1 \,.$$

for *Plancherel's theorem*, the square of the magnitude of Fourier transform $G(f)$ is unitary as well,

$$\int_{f=-\infty}^{+\infty} |G(f)|^2\, df = \int_{f=-\infty}^{+\infty} G^*(f)\, G(f)\, df = 1 \,,$$

and thus it can be interpreted as a probability density function in frequency domain. The following uncertainty relation holds

$$\sigma_{t,g}^2 \sigma_{f,G}^2 \geq \left( \frac{1}{2\pi} \frac{1}{2} \right)^2 \,,$$

or in terms of pulsation $\omega = 2\pi f$,

$$\sigma_{t,g}^2 \sigma_{\omega,G}^2 \geq \left( \frac{1}{2} \right)^2 \,,$$

**Heisenberg uncertainty relation in quantum mechanics**

Space and momentum representation of the state function $\Psi$ are related by the transformation,

$$\langle x|\Psi \rangle := \psi(x,t) = \int_{p=-\infty}^{+\infty} \psi_p(p,t)\, e^{i\frac{p}{\hbar}x}\, dp \,,$$

as it's shown in the section Quantum Mechanics:From position to momentum representation. The wave number reads $k = \frac{p}{\hbar}$. Starting from the uncertainty relation between the space coordinate $x$ and the wave number $k$,

$$\sigma_x \sigma_k \geq \frac{1}{2} \,,$$

Heisenberg uncertainty principle for position and momentum (for the same Cartesian coordinates) reads

$$\sigma_x \sigma_p \geq \frac{\hbar}{2} \ .$$

## Proof of the uncertainty relation

Assuming zero average $\bar{t} = 0$, $\bar{f} = 0$ (see below for proof without this assumption)

$$\sigma_{t,g}^2 \sigma_{f,G}^2 = \int_{t=-\infty}^{+\infty} |t|^2 \ g^*(t) \, g(t) \, dt \ \int_{f=-\infty}^{+\infty} |f|^2 \ G^*(f) \, G(f) \, df =$$

$$= \int_{t=-\infty}^{+\infty} |t \, g(t)|^2 \ dt \int_{f=-\infty}^{+\infty} |f \, G(f)|^2 \ df = \qquad (1)$$

$$= \int_{t=-\infty}^{+\infty} |t \, g(t)|^2 \ dt \int_{t=-\infty}^{+\infty} \left| -\frac{i}{2\pi} \, \dot{g}(t) \right|^2 \ dt \geq \qquad (2)$$

$$= \left| \int_{t=-\infty}^{+\infty} -t \, g^*(t) \frac{i}{2\pi} \, \dot{g}(t) \, dt \right|^2 \geq \qquad (3)$$

$$= \left( \frac{1}{2\pi} \right)^2 \left( \frac{1}{2} \right)^2$$

having used in

**(1)**

$$\mathcal{F}\{\dot{g}(t)\}(f) = \int_{t=-\infty}^{+\infty} \dot{g}(t) e^{-i2\pi ft} \, dt = \cdots = i2\pi \, f \, G(f) \ ,$$

$$f \, G(f) = -i \, \frac{\mathcal{F}\{\dot{g}(t)\}}{2\pi}$$

and thus *Plancherel's theorem*

$$\int_{f=-\infty}^{+\infty} |f \, G(f)|^2 \, df = \int_{f=-\infty}^{+\infty} \left| -\frac{i}{2\pi} \mathcal{F}\{\dot{g}(t)\} \right|^2 \, df =$$

$$= \int_{t=-\infty}^{+\infty} \left| -\frac{i}{2\pi} \, \dot{g}(t) \right|^2 \, dt$$

in (2) Cauchy-Schwartz inequality,

**(3)**

$$a := \int_{t=-\infty}^{+\infty} t g^*(t) \dot{g}(t) \, dt =$$

$$= \underbrace{[t \, g^*(t) \, g(t)] \, |_{-\infty}^{+\infty}}_{=0} - \int_{t=-\infty}^{+\infty} \frac{d}{dt} \left( t g^*(t) \right) g(t) \, dt =$$

$$= -\int_{t=-\infty}^{+\infty} g^*(t) \, g(t) \, dt - \int_{t=-\infty}^{+\infty} t \dot{g}^*(t) g(t) \, dt =$$

$$= -1 - a^* \ .$$

$$-1 = a + a^* = 2\,\mathrm{re}\{a\}\,,$$

and thus

$$|a|^2 \geq \mathrm{re}\{a\}^2 = \frac{1}{4}\,.$$

If $\bar{t} \neq 0$, or $\overline{f} \neq 0$,

$$
\begin{aligned}
\sigma_{t,g}^2 \sigma_{f,G}^2 &= \int_{t=-\infty}^{+\infty} \left|t - \bar{t}\right|^2 g^*(t)\,g(t)\,dt \int_{f=-\infty}^{+\infty} \left|f - \overline{f}\right|^2 G^*(f)\,G(f)\,df = \\
&= \int_{t=-\infty}^{+\infty} \left|(t - \bar{t})\,g(t)\right|^2 dt \int_{f=-\infty}^{+\infty} \left|(f - \overline{f})\,G(f)\right|^2 df = \qquad (1) \\
&= ...
\end{aligned}
$$

## 19.3 Relations between Fourier transforms

Some freestyle in changing order of summations and integrals, and use of generalized functions here…check it!

Different Fourier transforms exist, depending if the original function is:

- time discrete/time continuous
- periodic/non-periodic

namely,

- FS, Fourier series: time continuous, periodic function (or finite domain, with a periodic extension)
- FT, Fourier transform: time continuous, non-periodic function
- DTFT, discrete-time Fourier transform: time didscrete, infinite-length sequence
- DFT, discrete Fourier transform: time discrete, finite-length sequence (and then with a periodic extension)

### 19.3.1 Fourier transform of integrable functions

$$F(\nu) := \mathcal{F}\left\{f(t)\right\}(\nu) := \int_{t=-\infty}^{+\infty} f(t)e^{-i2\pi\nu t}\,dt\,,$$

### 19.3.2 Fourier transform of the sum of shifted integrable functions

The infinite sum of a shifted integrable function is defined as

$$\tilde{f}_T(t) = \sum_{n=-\infty}^{+\infty} f(t - nT)\,.$$

Its Fourier transform reads

$$\mathcal{F}\left\{\tilde{f}_T(t)\right\}(\nu) = \int_{t=-\infty}^{+\infty} \tilde{f}_T(t)e^{-i2\pi\nu t}\,dt =$$

$$= \sum_{n=-\infty}^{+\infty} \int_{t=-\infty}^{+\infty} f(t-nT)e^{-i2\pi\nu t}\,dt = \quad (1)$$

$$= \sum_{n=-\infty}^{+\infty} F(\nu)e^{-i2\pi\nu nT} = \quad (2)$$

$$= F(\nu) \sum_{n=-\infty}^{+\infty} e^{-i2\pi\nu nT} = \quad (3)$$

$$= \Delta\nu\, F(\nu)\, \mathrm{III}_{\Delta\nu}(\nu)\,,$$

having used properties of Fourier transform of shifted function in (1), and the properties of Dirac's comb in (3), having defined the frequency resolution

$$\Delta\nu := \frac{1}{T}\,.$$

This Fourier transform is proportional to the Fourier transform of the original function, sampled in frequency with elementary frequency $\Delta\nu$.

### 19.3.3 Fourier transform of the a function sampled with a Dirac comb - DTFT

Fourier transform of the original function sampled with $\Delta t\,\mathrm{III}_{\Delta t}(t)$ reads

$$\mathcal{F}\left\{\Delta t\, f(t)\,\mathrm{III}_{\Delta t}(t)\right\} = \Delta t \int_{t=-\infty}^{+\infty} f(t)\,\mathrm{III}_{\Delta t}(t)e^{-i2\pi\nu t}\,dt =$$

$$\sim \Delta t\frac{1}{\Delta t} \int_{t=-\infty}^{+\infty} f(t) \sum_{n=-\infty}^{+\infty} e^{in\frac{2\pi}{\Delta t}t}e^{-i2\pi\nu t}\,dt =$$

$$= \sum_{n=-\infty}^{+\infty} \int_{t=-\infty}^{+\infty} f(t)\, e^{-i2\pi(\nu-n\bar{\nu})t}\,dt = \quad (19.6)$$

$$= \sum_{n=-\infty}^{+\infty} F(\nu-n\bar{\nu}) = \mathrm{DTFT}(f(t);\Delta t)\,,$$

i.e. equals the periodic sum of the Fourier of the original function, with period

$$\bar{\nu} := \frac{1}{\Delta t}\,.$$

From this last sentence and from the *symmetry properties of Fourier transform*, **Nyquist-Shannon sampling theorem** follows seamlessly.

---

**Theorem 19.3.1 (Nyquist-Shannon sampling theorem)**

In order to **avoid aliasing** the sampling frequency must be twice the maxiumum[1] frequency in the signal,

$$\nu_s \geq 2\nu_{max}\,.$$

---

[1] Usually there's no such a frequency above which the signal is exactly zero, but usually there's a frequency above which the spectrum of the signal is approximately zero, i.e. below a threshold where it can be treated as zero, and introduce no aliasing.

---

**todo** check alternative expressions if using the definition of train of impulses instead of the Fourier series of Dirac's comb.

$$= \Delta t \int_{t=-\infty}^{+\infty} f(t) \sum_{k=-\infty}^{+\infty} \delta(t - k\Delta t) \, e^{-i2\pi\nu t} \, dt =$$

$$= \Delta t \sum_{k=-\infty}^{+\infty} f(k\Delta t) e^{-i2\pi\nu k\Delta t} = \text{DTFT}\left(f(t); \Delta t\right)$$

(19.7)

### 19.3.4 Fourier transform of the sum of shifted integral functions sampled with a Dirac comb

Fourier transform of the periodic sum

$$\Delta t \, \tilde{f}(t) \, \text{III}_{\Delta t}(t) = \Delta t \sum_{n=-\infty}^{+\infty} f(t - nT) \, \text{III}_{\Delta t}(t)$$

reads

$$\mathcal{F}\left\{ \Delta t \, \tilde{f}(t) \, \text{III}_{\Delta t}(t) \right\}(\nu) = \Delta t \int_{t=-\infty}^{+\infty} \sum_{n=-\infty}^{+\infty} f(t - nT) \sum_{k=-\infty}^{+\infty} \delta(t - k\Delta t) \, e^{-i2\pi\nu t} \, dt =$$

$$= \Delta t \sum_{n=-\infty}^{+\infty} \sum_{k=-\infty}^{+\infty} f(k\Delta t - nT) \, e^{-i2\pi\nu k\Delta t} =$$

and defining $k\Delta\tau_n := k\Delta t - nT$,

$$= \Delta t \sum_{n=-\infty}^{+\infty} \sum_{k=-\infty}^{+\infty} f(k\Delta\tau_n) e^{-i2\pi\nu k\Delta\tau_n} e^{-i2\pi\nu nT} =$$

$$= \Delta t \underbrace{\sum_{k=-\infty}^{+\infty} f(k\Delta\tau_n) e^{-i2\pi\nu k\Delta\tau_n}}_{=\text{DTFT}(f(t),\Delta t)} \underbrace{\sum_{n=-\infty}^{+\infty} e^{-i2\pi\nu nT}}_{=\Delta\nu \, \text{III}_{\Delta\nu}(\nu)} =$$

$$= \text{DTFT}(f(t), \Delta t) \, \Delta\nu \, \text{III}_{\Delta\nu}(\nu) \, .$$

**todo check!** check the change of coordinates that makes DTFT appear

**todo check!** what follows or, using the relation between $\Delta t$ and $T = N\Delta t$, $\Delta\nu = \frac{1}{T}$, and thus

$$\Delta t \, \Delta\nu = \Delta t \, \frac{1}{T} = \frac{1}{N} \, ,$$

it follows

$$= \frac{1}{N} \sum_{k=-\infty}^{+\infty} f(k\Delta\tau_n) e^{-i2\pi\nu k\Delta\tau_n} \, \text{III}_{\Delta\nu}(\nu) \, .$$

### 19.3.5 Useful properties

**Dirac's comb $\text{III}_T(t)$**

Dirac comb $\text{III}_T(t)$ is defined as a train of Dirac's delta

$$\text{III}_T(t) = \sum_{m=-\infty}^{+\infty} \delta(t - mT) \, .$$

Coefficients (19.2) of the Fourier series (19.1) of a $T$-periodic train of Dirac delta for $t \in \left[-\frac{T}{2}, \frac{T}{2}\right]$, read

$$c_n = \frac{1}{T} \int_{t=0}^{T} \delta(t) \, e^{-in\frac{2\pi}{T}t} = \frac{1}{T} \, ,$$

and thus the Fourier series of Dirac comb $\text{III}_T(t)$ reads

$$\text{III}_T(t) = \sum_{m=-\infty}^{+\infty} \delta(t - mT) \sim \frac{1}{T} \sum_{n=-\infty}^{+\infty} e^{in\frac{2\pi}{T}t} \, .$$

**Symmetry of Fourier transform**

# Part VII

# Calculus of Variations

# INTRODUCTION TO CALCULUS OF VARIATIONS

Calculus of variation deals with variations - i.e. "small changes" - of functions and functionals.

The meaning of the term functional may vary on the subfield of interest. In the field of calculus of variation, a **functional** can be defined as a function of function, i.e. a function whose argument is another function.

## Fields and applications

Fields and applications related to calculus of variations (give some examples below):

- gradient-based techniques like some methods in:
  - optimization, either free or constrained (via Lagrange multiplier methods)
  - sensitivity
- classical mechanics and physics in general:
  - analytical mechanics: Lagrangian formulation and Hamiltonian formulation of classical mechanics
- …

## Examples

- Lagrange equations for general problem
- examples:
  - brachistochrone for minimum time,…
  - catenary, i.e. static solution of wire and cables with neglibile bending stiffness
  - isoperimetric inequality, i.e. circle is the plane closed curve with given perimeter enclosing the largest area
- sensitivity of results to parameters. Some interesting sensitivity, both in time and trasnformed domains
  - characteristics of a system:
    - * equilibria
    - * eigenvalues
    - * …
- optimal control methods

## 20.1 Lagrange equations

Given the functional $S$, with arguments a function $q(t)$ and the independent variable $t$,

$$S[q(t), t] = \int_{t=t_0}^{t_1} L(\dot{q}(t),\, q(t),\, t)\, dt$$

its variation w.r.t. the function $q(t)$ reads

$$\delta S[q(t), t] = \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( S[q(t) + \varepsilon w(t),\, t] - S[q(t),\, t] \right)$$

where the function $w(t)$ is arbitrary, among those satisfying the constraint of the problems: as an example here, if the function $q(t)$ has prescribed values $q^*$ for some values of the independent variable, $t^*$, the variation $w(t)$ of the function $q(t)$ is zero there, $w(t^*)$ so that the variated function $q(t) + \varepsilon w(t)$ satisfies the constraint as well, i.e. $q(t^*) + \varepsilon w(t^*) = q^*$.

**Variation involves only small changes of function arguments**, since these ones are the elements that can be effectively changed, while the independent variable is not.

Direct computation of the variation gives

$$\delta S[q(t), t] = \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( S[q(t) + \varepsilon w(t),\, t] - S[q(t),\, t] \right) =$$

$$= \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( \int_{t=t_0}^{t_1} L(\dot{q}(t) + \varepsilon \dot{w}(t),\, q(t) + \varepsilon w(t),\, t) - \int_{t=t_0}^{t_1} L(\dot{q}(t),\, q(t),\, t)\, dt \right) =$$

$$= \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \int_{t=t_0}^{t_1} \left( L(\dot{q}(t) + \varepsilon \dot{w}(t),\, q(t) + \varepsilon w(t),\, t) - L(\dot{q}(t),\, q(t),\, t) \right)\, dt =$$

$$= \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \int_{t=t_0}^{t_1} \left\{ L(\dot{q}(t),\, q(t),\, t) + \varepsilon \left[ \frac{\partial L}{\partial \dot{q}} \dot{w}(t) + \frac{\partial L}{\partial q} w(t) \right] + o(\varepsilon) - L(\dot{q}(t),\, q(t),\, t) \right\}\, dt =$$

$$= \lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \int_{t=t_0}^{t_1} \left\{ \varepsilon \left[ \frac{\partial L}{\partial \dot{q}} \dot{w}(t) + \frac{\partial L}{\partial q} w(t) \right] + o(\varepsilon) \right\}\, dt =$$

$$= \int_{t=t_0}^{t_1} \left\{ \frac{\partial L}{\partial \dot{q}} \dot{w}(t) + \frac{\partial L}{\partial q} w(t) \right\}\, dt =$$

$$= \left[ w(t) \frac{\partial L}{\partial \dot{q}} \right]\Bigg|_{t=t_0}^{t_1} + \int_{t=t_0}^{t_1} \left\{ -\frac{d}{dt}\left( \frac{\partial L}{\partial \dot{q}} \right) + \frac{\partial L}{\partial q} \right\} w(t)\, dt \ .$$

The solution depends on the boundary conditions at the extreme points $t_0,\, t_1$. **If** the value of the function $q(t)$ is prescribed in $t_0$ and $t_1$, $q(t_0) = q_0,\, q(t_1) = q_1$, then its variation is zero, $w(t_0) = w(t_1) = 0$, for the reason that has been discussed above. The variation of the functional with prescribed boundary values of the argument function thus reads

$$\delta S[q(t), t] = \int_{t=t_0}^{t_1} \left\{ -\frac{d}{dt}\left( \frac{\partial L}{\partial \dot{q}} \right) + \frac{\partial L}{\partial q} \right\} \delta q(t)\, dt \ ,$$

having called $w(t) =: \delta q(t)$ to stress that is the variation of function $q(t)$. This notation - it's just notation, it has no special properties - could be useful if the functional depends on several arguments.

**Stationary conditions, $\delta S = 0$.** Stationary condition of the functional $S$ implies that $\delta S = 0$ for all the possible variations of the argument function, $\forall \delta q(t)$. This condition implies that the integrand is identically zero, i.e. **Lagrange equations**,

$$\frac{d}{dt}\left( \frac{\partial L}{\partial \dot{q}} \right) - \frac{\partial L}{\partial q} = 0 \ ,$$

### Higher-order derivatives

**Method 1.** If the Lagrangian function $L$ depends on higher order derivatives,

$$L\left(q^{(n)}(t),\, q^{(n-1)}(t),\, \dots,\, q'(t),\, q(t),\, t\right)$$

it's possible to recast the problem defining the $n$-dimensional function, $\mathbf{q}(t)$,

$$\mathbf{q}(t) = (q^0(t), q^1(t), \dots, q^{n-1}(t)) := \left(q(t), q'(t), \dots, q^{(n-1)}(t)\right) \ .$$

With some abuse of notation in $L$, the functional $S$ can be recasted as

$$S[q(t), t] = \int_{t=t_0}^{t_1} L(q^{(n)}(t),\, \dots,\, q(t),\, t)\, dt =$$

$$= \int_{t=t_0}^{t_1} L(\dot{\mathbf{q}}(t),\, \mathbf{q}(t),\, t)\, dt \ .$$

**todo** *Add constraints on components of $\mathbf{q}(t)$?*

Repeating the computation, the variation of the functional reads

$$\delta S[\mathbf{q}(t), t] = \left[\delta\mathbf{q}^T(t)\frac{\partial L}{\partial \dot{\mathbf{q}}}\right]\Bigg|_{t=t_0}^{t_1} + \int_{t=t_0}^{t_1} \delta\mathbf{q}^T(t) \left\{-\frac{d}{dt}\left(\frac{\partial L}{\partial \dot{\mathbf{q}}}\right) + \frac{\partial L}{\partial \mathbf{q}}\right\} dt \ .$$

**Method 2.** …

## 20.1.1 Euler-Beltrami equation

If the Lagrangian function $L$ is not an explicit function of the independent variable, $L(q'(x), q(x))$, Euler-Berltrami equation follows from the derivative of the Lagrangian,

$$\frac{dL}{dx} = \frac{\partial L}{\partial q'}q''(x) + \frac{\partial L}{\partial q}q' =$$

$$= \frac{\partial L}{\partial q'}q'' + \frac{d}{dx}\left(\frac{\partial L}{\partial q'}\right)q' =$$

$$= \frac{d}{dx}\left(\frac{\partial L}{\partial q'}q'\right) \ ,$$

and thus

$$\frac{d}{dx}\left[L - q'\frac{\partial L}{\partial q'}\right] = 0 \qquad \rightarrow \qquad L - q'\frac{\partial L}{\partial q'} = C \quad \text{const.}$$

**Note 1.** While Lagrange equations are a set of $N$ equations if the functional depends on $N$ argument functions $q_k(t)$, $k = 1 : N$, Euler-Beltrami equation is an equation only. Indeed for multiple argument functions

$$\frac{dL}{dx} = \frac{\partial L}{\partial q'_k}q''_k(x) + \frac{\partial L}{\partial q_k}q'_k =$$

$$= \frac{\partial L}{\partial q'_k}q''_k + \frac{d}{dx}\left(\frac{\partial L}{\partial q'_k}\right)q'_k = \frac{d}{dx}\left(\frac{\partial L}{\partial q'_k}q'_k\right) \ ,$$

where Einstein's summation notation of repeated index is used. Euler-Beltrami thus reads

$$L(q'_l(x), q_l(x)) - q'_k(x)\frac{\partial L}{\partial q'_k}(q'_l(x), q_l(x)) = C \ .$$

**Note 2.** If the Lagrangian function is an explicit function of the independent variable $x$, $L(q'(x), q(x), x)$, it's not hard to realize that the derivative of the Lagrangian function, along with the use of the Lagrange equation, gives

$$\frac{d}{dx}\left[L - q'\frac{\partial L}{\partial q'}\right] = \frac{\partial L}{\partial x} \ .$$

---

**Example 20.1.1 (Euler-Beltrami with $L(q'(x), q(x), x)$, Hamiltonian, energy and E.Noether)**

Euler-Beltrami equation shows that if $L(q'(x), q(x))$, thus $L - q'\partial_{q'}L$ is constant, (or an integral of motion in dynamics). In analytical mechanics (Lagrange mechanics, Hamiltonian mechanics), Lagrangian and Hamiltonian functions of a system read

$$L(\dot{q}_k(t), q_k(t), t) = T(\dot{q}, q, t) + U(q, t)$$

$$H(p, q, t) := p_k\dot{q}_k - L = \dot{q}_k\frac{\partial L}{\partial \dot{q}^k} - L \ ,$$

having used the common definition of the generalized momenta $p_k := \frac{\partial L}{\partial \dot{q}^k}$. It should be immediate to realize that the Hamiltonian is just the quantity appearing in Euler-Beltrami equation (or in its "modified version" if $\partial_t L \neq 0$), and thus

$$\frac{dH}{dt} = \frac{\partial L}{\partial t} \ .$$

In mechanics, if $\partial_t L = 0$, the Hamiltonian is a constant of motion. In this case, it can be prove that the Hamiltonian is equal to the eneergy of the system.

---

**Classical examples.**

---

**Example 20.1.2 (Brachistochrone)**

Find the trajectory…

- Elementary length: $ds = v\,dt$

- Energy: $E(y) = \frac{1}{2}mv^2 - mgy + C$. Setting $E = 0$ at starting point, from rest, at $y_0 = 0$, it implies $C = 0$; thus $v = \sqrt{2gy}$

- $x(s), y(s)$,

- $ds = \sqrt{dx^2 + dy^2} = \sqrt{1 + y'^2(x)}\,dx$

$$T = \int_{t_0}^{t_1} dt = \int_{s_0}^{s_1} \frac{ds}{v} = \int_{x_0}^{x_1} \frac{\sqrt{1 + y'^2(x)}}{\sqrt{2gy(x)}}\,dx$$

The Lagrangian doesn't explicitly depend on $x$, thus Euler-Beltrami equation can be used. Partial derivative of the Lagrangian function w.r.t. $q'$ reads

$$\frac{\partial L}{\partial y'} = \ldots = \frac{1}{\sqrt{2gy}\sqrt{1 + y'^2}}y' \ ,$$

and thus Euler-Beltrami equation reads

$$C = L - q'\frac{\partial L}{\partial y'} = \frac{1 + y'^2 - y'^2}{\sqrt{2gy}\sqrt{1 + y'^2}} = \frac{1}{\sqrt{2gy}\sqrt{1 + y'^2}}$$

Squaring $2gC^2 = \frac{1}{y(1+y'^2)}$, it's possible to write

$$y(x) = \frac{1}{2gC^2(1 + y'^2(x))} \ ,$$

Making the substitution $y(x)' = ...$

**Example 20.1.3 (Catenary)**

**Example 20.1.4 (Isoperimetric problem)**

# Part VIII

# Ordinary Differential Equations

# INTRODUCTION TO ORDINARY DIFFERENTIAL EQUATIONS

# LINEAR TIME-INVARIANT SYSTEMS

A linear time invariant system is governed by a linear ODE with constant coefficients. These equations can be recast as a first order system of ODEs,

$$
\begin{cases}
\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \\
\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{u} \\
\mathbf{x}(0^-) = \mathbf{x}_0
\end{cases}
$$

Exploiting the *properties of matrix exponential* the general expression of the state can be written as the sum of the free response to initial condition and the forced response.

$$
\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)\,d\tau
$$

$$
\mathbf{y}(t) = \mathbf{C}e^{\mathbf{A}t}\mathbf{x}_0 + \mathbf{C}\int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)\,d\tau + \mathbf{D}\mathbf{u}(t)
$$

### Proof in time domain

Multipying by $e^{-\mathbf{A}t}$,

$$
e^{-\mathbf{A}t}(\dot{x}(t) - \mathbf{A}\mathbf{x}(t)) = e^{-\mathbf{A}t}\mathbf{B}\mathbf{u}(t)
$$

$$
\frac{d}{dt}\left(\mathbf{x}e^{-\mathbf{A}t}\right) = e^{-\mathbf{A}t}\mathbf{B}\mathbf{u}(t)
$$

and integrating from $0^-$ to a generic time value $t$,

$$
e^{-\mathbf{A}t}\mathbf{x}(t) - \mathbf{x}_0 = \int_{\tau=0^-}^{t} e^{-\mathbf{A}\tau}\mathbf{B}\mathbf{u}(\tau)\,d\tau
$$

The state $\mathbf{x}(t)$ can be written as the sum of the free response and a force response. The general expression of the state and the output as a function reads

$$
\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)\,d\tau
$$

$$
\mathbf{y}(t) = \mathbf{C}e^{\mathbf{A}t}\mathbf{x}_0 + \mathbf{C}\int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)\,d\tau + \mathbf{D}\mathbf{u}(t)
$$

**Laplace domain.** The *Laplace transform* of the problem reads

$$
\begin{cases}
s\hat{\mathbf{x}} = \mathbf{A}\hat{\mathbf{x}} + \mathbf{B}\hat{\mathbf{u}} + \mathbf{x}_0 \\
\hat{\mathbf{y}} = \mathbf{C}\hat{\mathbf{x}} + \mathbf{D}\hat{\mathbf{u}}
\end{cases}
$$

$$(s\mathbf{I} - \mathbf{A})\hat{\mathbf{x}} = \mathbf{B}\hat{\mathbf{u}} + \mathbf{x}_0$$

$$\hat{\mathbf{x}}(s) = (s\mathbf{I} - \mathbf{A})^{-1}\mathbf{x}_0 + (s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}\hat{\mathbf{u}}(s)$$

$$\hat{\mathbf{y}}(s) = \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{x}_0 + \left[\mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}\right]\hat{\mathbf{u}}(s)$$

Performing inverse Laplace transform allows to go back to time domain (just use Laplace inverse transform of a matrix exponential, and the formula (18.1) for Laplace transform of convolution).

## 22.1 Impulsive force

The effect of an impulsive force $\mathbf{u}_\delta\delta(t)$ at time $t = 0$ is equivalent to an instantaneous change in the initial state, $\Delta\mathbf{x}_0 = \mathbf{B}\mathbf{u}_\delta$, from time $0^-$ before the impulse to time $0^+$ after the impulse. Splitting the input $\mathbf{u}(t)$ as the sum of impulsive input and regular input,

$$\mathbf{u}(t) = \mathbf{u}_r(t) + \mathbf{u}_\delta\delta(t)$$

$$\hat{\mathbf{u}}(s) = \hat{\mathbf{u}}_r(s) + \mathbf{u}_\delta$$

the solution in **time domain** reads

$$\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\left(\mathbf{u}_r(\tau) + \mathbf{u}_\delta\delta(\tau)\right)\,d\tau =$$

$$= e^{\mathbf{A}t}\left(\mathbf{x}_0 + \mathbf{B}\mathbf{u}_\delta\right) + \int_{\tau=0^-}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}_r(\tau)\,d\tau$$

$$\mathbf{y}(t) = \mathbf{C}e^{\mathbf{A}t}\left(\mathbf{x}_0 + \mathbf{B}\mathbf{u}_\delta\right) + \int_{\tau=0^-}^{t} \mathbf{C}e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}_r(\tau)\,d\tau + \mathbf{D}\mathbf{u}_r(t) + \mathbf{D}\mathbf{u}_\delta(t)$$

while in **Laplace domain** reads

$$\hat{\mathbf{x}}(s) = (s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}\hat{\mathbf{u}}_r(s) + (s\mathbf{I} - \mathbf{A})^{-1}\left(\mathbf{x}_0 + \mathbf{B}\mathbf{u}_\delta\right)$$

$$\hat{\mathbf{y}}(s) = \left[\mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}\right]\hat{\mathbf{u}}_r(s) + \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\left(\mathbf{x}_0 + \mathbf{B}\mathbf{u}_\delta\right) + \mathbf{D}\mathbf{u}_\delta$$

## 22.2 Properties

**Matrix exponential.**

$$e^{\mathbf{A}t} = \sum_{k=0}^{+\infty} \frac{\mathbf{A}^k t^k}{k!}\ .$$

Assuming it's possible swap derivative operator and summation (when?), it's possible to write

$$\frac{d}{dt}e^{\mathbf{A}t} = \frac{d}{dt}\sum_{k=0}^{+\infty} \frac{\mathbf{A}^k t^k}{k!} = \sum_{k=1}^{+\infty} kt^{k-1}\frac{\mathbf{A}^k t^{k-1}}{k!} = \mathbf{A}e^{\mathbf{A}t}\ .$$

**Laplace transform of exponential matrix.**

$$\mathcal{L}\left\{e^{\mathbf{A}t}\right\}(s) := \int_{t=0^-}^{+\infty} e^{\mathbf{A}t}e^{-st}\,dt =$$

$$= \int_{t=0^-}^{+\infty} e^{(-s\mathbf{I}+\mathbf{A})t}\,dt =$$

$$= (-s\mathbf{I} + \mathbf{A})^{-1}\ e^{(-s\mathbf{I}+\mathbf{A})t}\Big|_{t=0^-}^{+\infty} =$$

$$= (s\mathbf{I} - \mathbf{A})^{-1}\ ,$$

for all the values of $s$ for which $-s\mathbf{I} + \mathbf{A}$ is asymptotically stable, i.e. has all the eignevalues (thus, assuming that the matrix $\mathbf{A}$ can be diagonalizable. What happens if not? Exploit other matrix decompositions to draw conclusions) with negative real parts, and thus for all the values of $s > \max \mathrm{re}\{s_k(\mathbf{A})\}$, as it's shown in *Example 22.2.1*

---

**Example 22.2.1 (Asymptotic stability of a matrix A)**

An $N \times N$ diagnonalizable matrix $\mathbf{A}$,

$$\mathbf{A}\mathbf{v}_k = \mathbf{v}_k \, s_k \tag{22.1}$$

has all the eigenvalues with negative real part, $\mathrm{re}\{s_k\} < 0, \forall k = 1 : N$.

The eigenvalues of a matrix $a\mathbf{I} + \mathbf{A}$ are $a + s_k$, while the eigenvectors are the same as those of the matrix $\mathbf{A}$. This can be easily proved adding $a\mathbf{I}\mathbf{v}_k$ to both sides of equation (22.1),

$$(a\mathbf{I} + \mathbf{A}) \, \mathbf{v}_k = \mathbf{v}_k(a + s_k) \, .$$

---

**Transform of the convolution.**

# LTI SYSTEM RESPONSE

Usually the response of LTI to 3 different input are studied

- integrable input, being response to non-zero initial condition equivalent to an impulsive force at time $t = 0$, see *equivalence impulsive force - istantanteous change of state*
- periodic input
- stochastic input

## 23.1 LTI

$$\begin{cases} \dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu} \\ \mathbf{y} = \mathbf{Cx} + \mathbf{Du} \end{cases}$$

with initial conditions, $x(0) = \mathbf{x}_0$.

## 23.2 Response of LTI systems

The response of a LTI can be written as the sum of a free and forced response.

$$\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0}^{t} e^{\mathbf{A}(t-\tau)}\mathbf{Bu}(\tau)\, d\tau$$

$$\mathbf{y}(t) = \mathbf{C}\, e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0}^{t} \mathbf{C}e^{\mathbf{A}(t-\tau)}\mathbf{Bu}(\tau)\, d\tau + \mathbf{Du}(t)$$

(23.1)

**Proof**

**todo**

---

**Impulse response, $\mathbf{H}(t)$**

Functions

$$\mathbf{H}_x(t) := e^{\mathbf{A}t}\mathbf{B}$$

$$\mathbf{H}(t) := \mathbf{C}e^{\mathbf{A}t}\mathbf{B} + \mathbf{D}\delta(t)$$

are defined **impulse response functions** for the state $\mathbf{x}$ and the output $\mathbf{y}$, as they determine with the evolution of the state and the output of a system with zero initial condition $\mathbf{x}_0 = \mathbf{0}$ as a consequence of a impulsive force $\mathbf{u}_\delta(t) = \mathbf{u}\delta(t)$. Using Lagrange formula (23.1)

$$
\begin{aligned}
\mathbf{x}(t) &= e^{\mathbf{A}t}\mathbf{B}\mathbf{u} & &= \mathbf{H}_x(t)\mathbf{u} \\
\mathbf{y}(t) &= \left[\mathbf{C}e^{\mathbf{A}t}\mathbf{B} + \mathbf{D}\delta(t)\right]\mathbf{u} & &= \mathbf{H}(t)\mathbf{u}
\end{aligned}
\tag{23.2}
$$

**Convolution**

Integral in Lagrange formula (23.1) is the convolution of the impulse response function and the input. Changing the integration variable $z = t - \tau$, $d\tau = -dz$, $\tau = 0 : z = t$, $\tau = t$, $z = 0$ (and going back from $z$ to $\tau$ as the symbol for the dummy integration variable), it's possible to swap the arguments of the impulsive force and the forcing,

$$
\begin{aligned}
\mathbf{x}(t) &= e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0}^{t} \mathbf{H}_x(\tau)\mathbf{u}(t-\tau)\, d\tau \\
\mathbf{y}(t) &= \mathbf{C}\, e^{\mathbf{A}t}\mathbf{x}_0 + \int_{\tau=0}^{t} \mathbf{H}(\tau)\mathbf{u}(t-\tau)\, d\tau
\end{aligned}
\tag{23.3}
$$

**Linear but not time-invariant**

**todo** *The response of linear system, but not time-invariant can be written with the same expression used by Lagrange fomrula* (23.1), *with* $\mathbf{H}(t)$ *not equal to ..., but solution of the ODE...*

$$
\begin{aligned}
\Phi\left(\dot{\mathbf{x}} - \mathbf{A}\mathbf{x}\right) &= \Phi\mathbf{B}\mathbf{u} \\
\frac{d}{dt}\left(\Phi\mathbf{x}\right) &= \Phi\mathbf{B}\mathbf{u} \; ,
\end{aligned}
$$

with $\dot{\Phi} = -\Phi\mathbf{A}$, and proper initial conditions $\Phi(t, 0) = \mathbf{I}$.

...

$$
\mathbf{x}(t) = \Phi(t,0)\mathbf{x}_0 + \int_{\tau=0}^{t} \Phi(t,\tau)\mathbf{B}(\tau)\mathbf{u}(\tau)\, d\tau \; .
$$

## 23.3 Equilibria

Equilibrium solution of a ODE is a stationary solution, independent from time, $\dot{\mathbf{x}} = \mathbf{0}$. With a steady forcing $\overline{\mathbf{u}}$,

$$
\begin{cases}
\mathbf{0} = \mathbf{A}\overline{\mathbf{x}} + \mathbf{B}\overline{\mathbf{u}} \\
\mathbf{y} = \mathbf{C}\overline{\mathbf{x}} + \mathbf{D}\overline{\mathbf{u}} \; .
\end{cases}
$$

**If A is not singular**, for every steady forcing $\overline{\mathbf{u}}$ only one equilibrium exists for the system, $\overline{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{B}\overline{\mathbf{u}}$.

**If A is singular**...**todo** *rely on* Linear algebra

## 23.4 Stability

**Asymptotic stability to initial perturbations.** A system is stable under perturbation of initial conditions if the free response goes to zero as $t \to +\infty$,

$$\lim_{t \to +\infty} \mathbf{x}_{free}(t) = \mathbf{0} \ .$$

**todo** *for all perturbations*(?).

If matrix $\mathbf{A}$ is diagonalizable, a system is asymptotically stable if all the eigenvalues of the matrix $\mathbf{A}$ have negative real part, $\text{re}\{\lambda_i(\mathbf{A})\} < 0$.

**todo** *rely on* Linear algebra

## 23.5 Response to integrable input - and Laplace transform

$$\begin{aligned}
\mathbf{x}(s) &= (s\mathbf{I} - \mathbf{A})^{-1} \mathbf{x}_0 + \mathbf{H}_x(s)\mathbf{u}(s) \\
\mathbf{y}(s) &= \mathbf{C} \, (s\mathbf{I} - \mathbf{A})^{-1} \mathbf{x}_0 + \mathbf{H}(s)\mathbf{u}(s)
\end{aligned} \tag{23.4}$$

## 23.6 Response to periodic input of stable systems - and Fourier transform

Forced response to periodic input of an stable system can be represented in Fourier domain replacing $s$ with $j\omega$ (and dropping $j$ once it's clear we're dealing with Fourier analysis)

$$\begin{aligned}
\mathbf{x}(\omega) &= \mathbf{H}_x(\omega) \, \mathbf{u}(\omega) \\
\mathbf{y}(\omega) &= \mathbf{H}(\omega) \, \mathbf{u}(\omega)
\end{aligned} \tag{23.5}$$

## 23.7 Response to stochastic input of stable systems

### 23.7.1 Assumptions

…

### 23.7.2 Expected value

$$\mu_{\mathbf{x}}(t) := \mathbb{E}[\mathbf{x}(t)]$$

Taking expected value of LTI equations

$$\begin{cases}
\mathbb{E}[\dot{\mathbf{x}}] = \mathbf{A}\mathbb{E}[\mathbf{x}] + \mathbf{B}\mathbb{E}[\mathbf{u}] \\
\mathbb{E}[\mathbf{y}] = \mathbf{C}\mathbb{E}[\mathbf{x}] + \mathbf{D}\mathbb{E}[\mathbf{u}]
\end{cases}$$

**todo**

- *is it possible to swap time derivative and expected value operators?*

- *does time derivative of a stochastic process always exist? What's the right way to interpret time derivative in stochastic equations? As an increment?*

$$d\mathbf{x} = \mathbf{A}\mathbf{x}\,dt + \mathbf{B}\mathbf{u}\,dt$$

so that if $\mathbf{u}(t) \sim \xi(t)$, then $\mathbf{u}\,dt \sim d\mathbf{W}$

# 23.8 Analysis of response of LTI in Fourier domain

**Energy or power.** For

- time integrable signals: don't average in time; energy

- periodic or stochastic non-integrable signals: average in time; power

**Stationary stochastic process.** Strong-sense (joint probability independent from time) or wide-sense (moments up to order $m$ are independent on time)

**Ergodic process.** Average over random process can be replaced with average over time.

---

**Example of stationary non-ergodic process**

As an example of stationary, but not ergodic process is a random process whose realizations are

$$X_t = \theta\,, \qquad \forall t \in [0, T]$$

with $\theta$ drawn from a random variable $\Theta$, with average $\mu$. Average in time over a realization gives $\theta$, while average over realizations gives $\mu$ for all $t$,

$$\text{Avg. in time over a realization:} \quad \frac{1}{T}\int_{t=0}^{T} X_t\,dt = \frac{1}{T}\int_{t=0}^{T} \theta\,dt = \theta$$
$$\text{Avg. in probability over realizations:} \quad \mathbb{E}[X_t] = \mathbb{E}[\Theta] = \mu\,.$$

---

## 23.8.1 Stationary Ergodic processes

**Auto-correlation**

$$\mathbf{R}_{XX}(t_1, t_2) = \mathbb{E}\left[\mathbf{X}_{t_1}\,\mathbf{X}_{t_2}^*\right]$$

**Auto-covariance**

$$\mathbf{K}_{XX}(t_1, t_2) = \mathbb{E}\left[(\mathbf{X}_{t_1} - \mu_{t_1})(\mathbf{X}_{t_2} - \mu_{t_2})^*\right]$$

For stationary processes, both auto-correlation and auto-covariance don't depend on two time instant $t_1$, $t_2$ but only on the time-shift between them, $t_2 = t_1 + \tau$,

$$\mathbf{R}_{XX}(\tau) = \mathbb{E}\left[\mathbf{X}_t\,\mathbf{X}_{t+\tau}^*\right]\,, \qquad \forall t\,.$$

## Autocorrelation and ESD of time integrable processes

With causal, $\mathbf{x}(t \leq 0) = 0$, and integrable processes, autocorrelation is defined as

$$\tilde{\mathbf{R}}_{xx}(\tau) = \int_{t=0}^{+\infty} \mathbf{x}(t)\mathbf{x}^*(t+\tau)\,dt\ ,$$

**Energy Spectral Density (ESD)** is defined as its Fourier transform

$$\tilde{\Phi}_{xx}(f) = \mathcal{F}\left\{\mathbf{R}_{xx}(\tau)\right\}(f)\ ,$$

## Autocorrelation and PSD of infinite-time processes

Using average over time as an approximation of the probability average

$$\mathbb{E}[f(t)] \sim \lim_{T \to +\infty} \frac{1}{T} \int_{t=0}^{T} f(\tau)\,d\tau\ ,$$

**Power Spectral Density (PSD)**, defined as the Fourier transform of the auto-correlation, reads

$$\Phi_{YY}(f) = \mathcal{F}[\mathbf{R}_{YY}(\tau)](f) =$$
$$= \int_{\tau=-\infty}^{+\infty} \mathbb{E}\left[\mathbf{y}(t)\mathbf{y}^*(t+\tau)\right] e^{-i2\pi f\tau}\,d\tau\ .$$

For asymptotically stable systems subject to periodic forcing, response to non-homogenoeus initial conditions vanishes while the forced response can be written either in Fourier or in time domain, exploiting the property of the transform of convolution *Property 19.2.1*, as

$$\mathbf{y}(f) = \mathbf{H}(f)\,\mathbf{u}(f) \qquad , \qquad \mathbf{y}(t) = \mathbf{H} * \mathbf{u}(t) = \int_{\tau=-\infty}^{+\infty} \mathbf{H}(\tau)\mathbf{u}(t-\tau)\,d\tau\ .$$

Inserting harmonic response into the expression of PSD, it becomes

$$\Phi_{YY}(f) = \int_{\tau=-\infty}^{+\infty} \mathbb{E}\left[\left(\int_{\xi=-\infty}^{+\infty} \mathbf{H}(\xi)\,\mathbf{u}(t-\xi)\,d\xi\right)\left(\int_{\chi=-\infty}^{+\infty} \mathbf{H}(\chi)\,\mathbf{u}(t+\tau-\chi)\,d\chi\right)^*\right] e^{-i2\pi f\tau}\,d\tau = \quad (1)$$

$$= \int_{\tau=-\infty}^{+\infty}\int_{\xi=-\infty}^{+\infty}\int_{\chi=-\infty}^{+\infty} \mathbf{H}(\xi)\mathbb{E}\left[\mathbf{u}(t-\xi)\mathbf{u}^*(t+\tau-\chi)\right]\mathbf{H}^*(\chi)\,e^{-i2\pi f\tau}\,d\tau\,d\chi\,d\xi = \quad (2)$$

$$= \int_{\tau=-\infty}^{+\infty}\int_{\xi=-\infty}^{+\infty}\int_{\chi=-\infty}^{+\infty} \mathbf{H}(\xi)\mathbb{E}\left[\mathbf{u}(t)\mathbf{u}^*(t+\tau-\chi+\xi)\right]\mathbf{H}^*(\chi)\,e^{-i2\pi f\tau}\,d\tau\,d\chi\,d\xi = \quad (3)$$

$$= \int_{z=-\infty}^{+\infty}\int_{\xi=-\infty}^{+\infty}\int_{\chi=-\infty}^{+\infty} \mathbf{H}(\xi)\mathbb{E}\left[\mathbf{u}(t)\mathbf{u}^*(t+z)\right]\mathbf{H}^*(\chi)\,e^{-i2\pi f(z-\xi+\chi)}\,dz\,d\chi\,d\xi =$$

$$= \int_{\xi=-\infty}^{+\infty} \mathbf{H}(\xi)\,e^{i2\pi f\xi}\,d\xi \int_{z=-\infty}^{+\infty} \mathbf{R}_{uu}(z)\,e^{-i2\pi fz}\,dz \int_{\chi=-\infty}^{+\infty} \mathbf{H}^*(\chi)\,e^{-i2\pi f\chi}\,d\chi = (4)$$

$$= \overline{\mathbf{H}}(f)\,\Phi_{uu}(f)\,\mathbf{H}^T(f)\ .$$

having used (1) the fact that the impulse response function is deterministic, (2) time shift in the autocorrelation for a **stationary** random process $\mathbf{u}(t)$ as an input, and (3) the change of coordinates $(z, \xi, \chi) = (\tau - \chi + \xi, \xi, \chi)$, and (4) a slight notation abuse ndicating Fourier transform with the same symbols of the functions in time domain.

---

**23.8. Analysis of response of LTI in Fourier domain**

# TWENTYFOUR

# LTI: STABILITY AND FEEDBACK

## Contents

- Stability of a SISO LTI, w.r.t. non-zero i.c. or impulsive input
- 
- 

## 24.1 Stability of a LTI - SISO

### 24.1.1 Transfer function

Let the SISO input-output relation of a LTI system be represented in Laplace domain by the transfer function $G(s)$,

$$y(s) = G(s)u(s)$$

For rational transfer functions

$$G(s) = k\frac{\prod_{n=1}^{N}(z_n - s)}{\prod_{d=1}^{D}(p_d - s)} = \frac{\sum_{n=1}^{N} a_n s^n}{\sum_{d=0}^{D} b_d s^d} \ ,$$

with $z_n$ the zeros, $p_d$ the poles, and $k = \frac{a_0}{b_0} = \lim_{s \to 0} G(s)$ the static gain. If the system is strictly proper, $N < D$, and the TF can be written as a **sum of partial functions**. As an example, for a transfer function with simple poles

$$G(s) = \sum_{d=1}^{N} \frac{A_d}{(p_d - s)} \ ,$$

while for a pole $p_d$ with multiplicity $m_d > 1$ all the terms $\propto \frac{1}{(p_d-s)^{e_d}}$, with $e_d = 1 : m_d$ must be included, see example below.

$$G(s) = \sum_{d=1}^{N} \frac{A_d}{(p_d - s)^{e_d}} \ ,$$

**Example 24.1.1 (Sum of partial fractions of a TF with poles with multiplicity $> 1$)**

Let's write the rational function $G(s) = \frac{s+1}{(s+2)^3}$ as a sum of partial functions,

$$G(s) = \frac{s+1}{(s+2)^3} =$$
$$= \frac{A}{s+2} + \frac{B}{(s+2)^2} + \frac{C}{(s+2)^3} =$$
$$= \frac{A(s+2)^2 + B(s+2) + C}{(s+2)^3} =$$
$$= \frac{s^2(A) + s(4A+B) + 4A + 2B + C}{(s+2)^3} ,$$

The value of the coefficients $A$, $B$, $C$ is computed comparing the first and the last expression of the numerator of $G(s)$

$$\begin{cases} s^2 : & A = 0 \\ s \ : & 4A + B = 1 \\ 1 \ : & 4A + 2B + C = 1 \end{cases} \qquad \rightarrow \qquad \begin{cases} A = 0 \\ B = 1 \\ C = -1 \end{cases}$$

and thus

$$G(s) = \frac{s+1}{(s+2)^3} = \frac{1}{(s+2)^2} - \frac{1}{(s+2)^3} .$$

### 24.1.2 Stability w.r.t. non-zero initial conditions - or w.r.t. implusive input - in time domain

Transfer function $G(s)$ represents the free the function w.r.t. impulsive input. Response in time domain can be evaluated as the inverse Laplace transform of the transfer function,

$$y(t) = \mathcal{L}^{-1}\{y(s)\} = \mathcal{L}^{-1}\left\{\sum_{n_d=1}^{D} \frac{R_d}{(s - p_d)}\right\} = \sum_{d=1}^{D} R_d\, e^{p_d t} .$$

If $\mathrm{re}\{p_d\} < 0$ for $\forall d$, then the response is asymptotically stble for $t \to +\infty$, as $|e^{p_d t}| = e^{\mathrm{re}\{p_d\}t} \to 0$, for $t \to +\infty$.

## 24.2 Stability of closed-loop systems

### 24.2.1 Cauchy argument principle

For the *Cauchy argument principle*, the difference of argument of function $F(s)$ when $s$ performs a counter-clockwise loop over th contour $\Gamma$ enclosing poles $p_n$, and zeros $z_d$ of the function $F(s)$ reads

$$\Delta\arg\{F(s)\} = \sum_n \arg(s - z_n) - \sum_d \arg(s - p_d) = Z2\pi - P2\pi = (Z - P)2\pi ,$$

and thus the the diagram of function $F(s)$ performs

$$N = Z - P \tag{24.1}$$

loops around the origin $0 + i0$ of the complex plane.

## 24.2.2 Nyquist stability criterion

Nyquist stability criterion provides some conditions for the stability of a closed loop transfer function

$$G_c(s) := \frac{G(s)}{1 + G(s)} \, ,$$

being $G(s)$ the open loop transfer function, and the feedback with the opposite of the output ($y = G(s)\,(u - y)$).

The poles of the closed loop systems are the zeros of the function $1 + G(s)$. If the closed-loop system is asymptotically stable, it must have no pole with positive real part. As the complex variable $s$ performs a clockwise (and thus, the signs of the relation change w.r.t. Cauchy argument criterion) loop over **Nyquist path** (semicircle in the RHS half-plane; then **todo** discuss the case where poles and zeros are on the imaginary axis…deform the path…). No zero of $1 + G(s)$ with positive real part means $Z = 0$. Using the relation between poles, zeros and loops around the origin given by Cauchy argument principle (24.1), and recalling the opposite direction, it follows that

$$N = P \, ,$$

in order to have $Z = 0$, i.e. the diagram of $1 + G(s)$ must perform $N$ counter-clockwise loops around $0 + i0$ equal to the number of its poles with positive real parts. As the poles of $1 + G(s)$ are the same as the poles of $G(s)$, it's possible to formulate Nyquist criterion looking at $G(s)$, as

---

**Theorem 24.2.1 (Nyquist criterion)**

In order for the closed-loop system to be asymptotically stable, the diagram of the open-loop transfer function $G(s)$ must perform a number $N$ of counter-clocwise loops around the *critical point* $-1 + i\,0$ equal to the number of its zeros with positve real part.

---

## 24.2.3 Bode stability criterion for minimal phase systems

---

**Definition 24.2.1 (Minimal phase systems)**

---

…

Safetry gain margin, safety phase margin…

# Part IX

# Partial Differential Equations

# INTRODUCTION TO PARTIAL DIFFERENTIAL EQUATIONS

Partial differential equations usually comes from balance equations in **continuum mechanics**. Integral equations are the most general form of these equations, and an equivalent differential problem only exists if the fields involved in the equations are regular enough, for their derivatives to exist - and to apply theorems requiring some regularity of the functions.

Classical numerical methods:

- **FVM**: directly solves the **integral problem**, solving integral balance equations for cells in which the domain is divided

- **FDM**: given the problem in **differential form**, FDM directly approximates space derivatives of the **strong formulation** of the problem

- **FEM**: given the problem in **differential form**, FEM projects the **weak formulation** of the problem on a finite-dimensional space

- **BEM**: *integro-differential equation*, *singularities*,…

- **Spectral methods**,…

- **SEM**,…

## 25.1 Examples

In Physics:

- Advection equation

$$\partial_t u + \vec{a} \cdot \nabla u = f$$

- Diffusion equation

$$\partial_t u - \nu \nabla^2 u = f$$

- Hyperbolic equation/system of equations

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{F}(\mathbf{u}) = \mathbf{f}$$

- Wave equation

$$\frac{1}{c^2} \partial_{tt} u - \nabla^2 u = f$$

## 25.2 Balance equations in physics

- Small-strain continuum mechanics

$$\rho \partial_{tt} \vec{s} = \rho_0 \vec{g} + \nabla \cdot \boldsymbol{\sigma}(\boldsymbol{\varepsilon})$$

- Heat conduction

- Fluid dynamics

  - Navier-Stokes for compressible fluids (conservative or convective equations)

  $$\{$$

  - Navier-Stokes for incompressible fluids (convective form,…)

  $$\begin{cases} \rho \partial_t \vec{u} + \rho(\vec{u} \cdot \nabla)\vec{u} - \mu \nabla^2 \vec{u} + \nabla P = \rho \vec{g} \\ \nabla \cdot \vec{u} = 0 \end{cases}$$

**todo**

- Different forms of equations may be more or less convenient for different solution approaches

- Most of the physical laws comes from integral balance equation of the form

$$\frac{d}{dt} \int_{V_t} \rho \mathbf{u} = \int_{V_t} \rho \mathbf{r} + \oint_{\partial V_t} \hat{n} \cdot \mathbf{T}(\mathbf{u})$$

whose local - differential - form (in case of differentiable functions) readily follows from the application of Reynolds' transport theorem and divergence theorem to transform time derivative and boundary terms

$$\partial_t (\rho \mathbf{u}) + \nabla \cdot (\rho \mathbf{u} \vec{u}) = \rho \mathbf{r} + \nabla \cdot \mathbf{T}(\mathbf{u})$$
$$\partial_t (\rho \mathbf{u}) + \nabla \cdot \mathbf{F}(\mathbf{u}) = \rho \mathbf{r} \,,$$

and the physical meaning of each term is evident and readily expalinable as flux or volume or surface sources.

- Further manipulation/simplification may cover the clear meaning of the terms of the differential equation. As an example, the conservative form of Navier-Stokes equations for incompressible fluids with constant and uniform density read

$$\begin{cases} \partial_t (\rho \vec{u}) + \nabla \cdot (\rho \vec{u} \otimes \vec{u}) = \rho \vec{g} + \nabla \cdot \mathbb{T} \\ \nabla \cdot \vec{u} = 0 \,, \end{cases}$$

where the stress tensor for a Newtonian fluid reads

$$\begin{aligned} \mathbb{T} &= -p \mathbb{I} + 2\mu \mathbb{D} + \lambda \left( \nabla \cdot \vec{u} \right) \mathbb{I} \\ &= -p \mathbb{I} + \mu \left( \nabla \vec{u} + \nabla^T \vec{u} \right) + \lambda \left( \nabla \cdot \vec{u} \right) \mathbb{I} \end{aligned}$$

Using the incompressibility constraint $\nabla \cdot \vec{u}$, and treating the density $\rho$ as a constant and uniform parameter, the convective form of the Navier-Stokes equations reads

$$\begin{cases} \rho \partial_t \vec{u} + \rho \left( \vec{u} \cdot \nabla \right) \vec{u} = \rho \vec{g} - \nabla P + 2\mu \nabla \cdot \mathbb{D} \\ \nabla \cdot \vec{u} = 0 \,. \end{cases}$$

The divergence of the viscous stress tensor becomes

$$2\mu \nabla \cdot \mathbb{D} = \mu \nabla \cdot \left( \nabla \vec{u} + \nabla^T \vec{u} \right) = \mu \left( \nabla^2 \vec{u} + \nabla \underbrace{\left( \nabla \cdot \vec{u} \right)}_{=0} \right) = \mu \nabla^2 \vec{u} \,,$$

so that one of the most common form of incompressible Navier-Stokes equations follows

$$\begin{cases} \rho \partial_t \vec{u} + \rho(\vec{u} \cdot \nabla)\vec{u} - \mu \nabla^2 \vec{u} + \nabla P = \rho \vec{g} \\ \nabla \cdot \vec{u} = 0 \ . \end{cases}$$

It should be evident that in the latter form of Navier-Stokes equations no divergence explicitly appears, so that the right expression of surface source terms can't be found immediately. In momentum equation, surface source terms come from surface stress acting on the boundary of the domain, whose expression reads

$$\begin{aligned} \vec{t}_n = \hat{n} \cdot \mathbb{T} &= \\ &= \hat{n} \cdot (-P\mathbb{I} + 2\mu\mathbb{D}) = \\ &= \hat{n} \cdot (-P\mathbb{I} + \mu\left(\nabla\vec{u} + \nabla^T\vec{u}\right)) = \\ &= -P\hat{n} + \hat{n} \cdot (\mu\left(\nabla\vec{u} + \nabla^T\vec{u}\right)) = \end{aligned}$$

As an example, in weak formulation of incompressible Navier-Stokes problem the **natural boundary condition** arising in the method depends on the expression of the strong formulation of the NS problem. If one needs to prescribe stress boundary conditions, it could be an idea to start from NS equations w/o extra simplifications.

# ELLIPTIC EQUATIONS

## 26.1 Poisson equation

Given the volume density source $f(\vec{r})$ and the diffusivity $\nu(\vec{r})$, Poisson equation for the scalar field $\phi(\vec{r})$ reads

$$-\nabla \cdot (\nu \nabla \phi) = f \qquad \vec{r} \in V$$

with proper boundary conditions on $\partial V$. As an example, tipical boundary conditions are:

$$
\begin{aligned}
\phi(\vec{r}) = g(\vec{r}) && \vec{r} \in S_D && \text{esserntial - Dirichlet b.c.} \\
\nu \hat{n} \cdot \nabla \phi(\vec{r}) = h(\vec{r}) && \vec{r} \in S_N && \text{natural - Neumann b.c.} \\
a\phi(\vec{r}) + \nu \hat{n} \cdot \nabla \phi(\vec{r}) = b(\vec{r}) && \vec{r} \in S_R && \text{Robin b.c.}
\end{aligned}
$$

### Numerical methods

1-dimensional Poisson equation:

- *Finite Element Methods*
- *Finite Volume Methods*
- *Finite Difference Methods*
- *Boundary Element Methods*
- *Spectral Methods* and *Spectral Element Methods*

## 26.1.1 Weak formulation

For $\forall w \in \dots$ (functional space, recall some results about existence and uniqueness of the solution, Lax-Milgram theorem,…)

$$
\begin{aligned}
0 &= \int_V w \left\{ \nabla \cdot (\nu \nabla \phi) + f \right\} = \\
&= \oint_{\partial V} w \hat{n} \cdot (\nu \nabla \phi) + \int_V \left\{ -\nu \nabla \vec{w} \cdot \nabla \phi + wf \right\} =
\end{aligned}
$$

Splitting boundary contribution as the sum from single contributions from different regions, and applying boundary conditions, setting $w = 0$ for $\vec{r} \in S_D$ (see the ways to prescribe essential boundary conditions),

$$0 = \int_{S_D} \underbrace{w}_{=0} \hat{n} \cdot (\nu \nabla \phi) + \int_{S_N} w \underbrace{\hat{n} \cdot (\nu \nabla \phi)}_{=h} + \int_{S_R} w \underbrace{\hat{n} \cdot (\nu \nabla \phi)}_{=b-a\phi} + \int_V \left\{ -\nu \nabla \vec{w} \cdot \nabla \phi + wf \right\} .$$

and rearranging the equation separating terms containing unknowns from known contributions,

$$\int_V \nu \nabla w \cdot \nabla \phi + \int_{S_R} wa\phi = \int_V wf + \int_{S_N} wh + \int_{S_R} wb \qquad \forall w \in \dots \ ,$$

and $\phi = g$, for $\vec{r} \in S_D$.

### Different ways to prescribe essential boundary conditions

**Strong formulation.**

**Using Lagrance multiplier - weak formulation of essential boundary conditions.** Adding a the essential boundary condition as a constraint with Lagrange multipliers in the weak formulation of the problem,

$$\dots + \int_{S_D} w_D(\phi - g) \ ,$$

…

## 26.1.2 Existence and uniqueness

Assuming two solutions $u_1(\mathbf{r})$, $u_2(\mathbf{r})$ exist for the Poisson problem

$$\begin{cases} -\nabla \cdot (\nu \nabla u) = f & \mathbf{r} \in V \\ u|_{S_D} = g \\ \nu \hat{n} \cdot \nabla u|_{S_N} = h \end{cases}$$

Their difference $\delta u(\mathbf{r}) := u_2(\mathbf{r}) - u_1(\mathbf{r})$ then satisfies the homogeneous problem

$$\begin{cases} -\nabla \cdot (\nu \nabla \delta u) = 0 & \mathbf{r} \in V \\ \delta u|_{S_D} = 0 \\ \nu \hat{n} \cdot \nabla \delta u|_{S_N} = 0 \end{cases}$$

The norm of the gradient of the difference of the solution reads

$$\begin{aligned} \int_V \nu |\nabla \delta u|^2 = \int_V \nu \nabla \delta u \cdot \nabla \delta u = \\ = \int_V \nabla \cdot (\delta u\, \nu \nabla \delta u) - \int_V \delta u\, \underbrace{\nabla \cdot (\nu \nabla \delta u)}_{=0} = \\ = \oint_{\partial V} \delta u\, \nu \hat{n} \cdot \nabla \delta u = \\ = \int_{S_D} \underbrace{\delta u}_{\delta u|_{S_D}=0} \nu \hat{n} \cdot \nabla \delta u + \int_{S_N} \delta u\, \underbrace{\nu \hat{n} \cdot \nabla \delta u}_{\nu \hat{n} \cdot \nabla \delta u|_{S_N}=0} = \\ = 0 \ . \end{aligned}$$

**If** $\nu(\mathbf{r}) > 0$ for $\forall \mathbf{r} \in V^1$, it follows that

$$\nabla \delta u = \nabla (u_2 - u_1) = 0 \ ,$$

---

[1] As an example, this condition appears in physical systems with non-zero and positive diffusion coefficients in diffusion problems, like thermal conduction or specie diffusion via Fick's law.

and thus the two solutions differs at most by an additive constant,

$$u_2(\mathbf{r}) - u_1(\mathbf{r}) = c .$$

**If** the Dirichlet boundary has non-null dimension, it forces the value of the functions to coincide on that boundary, $u_2(\mathbf{r}) = u_1(\mathbf{r})$ on $S_D$, and thus sets the value of the additive constant $c$ to be zero, $c = 0$, and

$$u_2(\mathbf{r}) = u_1(\mathbf{r}) ,$$

thus proving the uniqueness of the solution of the Poisson problem, with non-negative diffusion coefficient and non-zero dimension of Dirichlet boundary.

# PARABOLIC EQUATIONS

## 27.1 Heat equation

Heat equation for a scalar field $\phi(\vec{r}, t)$ can be interpreted as the unsteady equation of a *Poisson equation*,

$$\partial_t \phi - \nabla \cdot (\nu \nabla \phi) = f \qquad (\vec{r}, t) \in V \times [0, T] \,,$$

with proper boundary and initial conditions, $\phi(\vec{r}, 0) = \phi_0(\vec{r})$. Common boundary conditions are the same as the one discussed for Poisson problem.

### 27.1.1 Weak formulation

For $\forall w \in \dots$ (functional space, recall some results about existence and uniqueness of the solution, Lax-Milgram theorem,…)

$$0 = \int_V w \left\{ -\partial_t \phi + \nabla \cdot (\nu \nabla \phi) + f \right\} =$$
$$= \oint_{\partial V} w \hat{n} \cdot (\nu \nabla \phi) + \int_V \left\{ -\partial_t \phi - \nu \nabla \vec{w} \cdot \nabla \phi + w f \right\} =$$

Splitting boundary contribution as the sum from single contributions from different regions, and applying boundary conditions, setting $w = 0$ for $\vec{r} \in S_D$ (see the ways to prescribe essential boundary conditions),

$$0 = \int_{S_D = 0} \underaccent{\smile}{w} \, \hat{n} \cdot (\nu \nabla \phi) + \int_{S_N} w \underbrace{\hat{n} \cdot (\nu \nabla \phi)}_{=h} + \int_{S_R} w \underbrace{\hat{n} \cdot (\nu \nabla \phi)}_{=k-\phi} + \int_V \left\{ -\partial_t \phi - \nu \nabla \vec{w} \cdot \nabla \phi + w f \right\} \,.$$

and rearranging the equation separating terms containing unknowns from known contributions,

$$\int_V w \partial_t \phi + \int_V \nu \nabla w \cdot \nabla \phi + \int_{S_R} w \phi = \int_V w f + \int_{S_N} w h + \int_{S_R} w k \qquad \forall w \in \dots \,,$$

and $\phi = g$, for $\vec{r} \in S_D$.

# TWENTYEIGHT

# HYPERBOLIC PROBLEMS

Hyperbolic problems often come from a small-amplitude linearization, or as the non-diffusion (or inviscid) limit of a more general problem.

As a result of these simplification, these problems may experience **shocks** (i.e. discontinuity in the solution, where the differential equations stop to hold, and integral equations and jump conditions are required). **todo** *classification of discontinuities on the massflow across the surface*

The very nature of these problem also suggest methods for the solution or the analysis of these equations, like **characteristic method**.

## 28.1 Scalar linear

### 28.1.1 1-dimensional

$$\partial_t u(x,t) + a \partial_x u(x,t) = f(x,t)$$

**Caracteristic method.** $U(t) = u(X(t), t)$, with the caracteristic curves $X(t)$ defined as those curves where the PDE becomes a ODE. Evaluating the time derivative of the function $u(X(t), t)$, the hyperbolic equation can be recast as

$$\frac{dU}{dt} + \left[ a(X(t), t) - \frac{dX}{dt} \right] \partial_x u = f(X(t), t) .$$

The equation of characteristic lines is

$$\frac{dX}{dt} = a(X(t), t) ,$$

and the PDE on characteristic line becomes the ODE

$$\frac{dU}{dt}(X(t), t) = f(X(t), t) .$$

## 28.2 Scalar non-linear

## 28.3 System linear

### 28.3.1 1-dimensional

$$\mathbf{u}(x,t)$$

$$\partial_t \mathbf{u} + \mathbf{A} \partial_x \mathbf{u} = \mathbf{f}$$

## Method of characteristics

**Characteristics.** $\mathbf{U}(t) = \mathbf{u}(X(t), t)$

$$\frac{d\mathbf{U}}{dt} - \frac{dX}{dt}\partial_x\mathbf{u} + \mathbf{A}\partial_x\mathbf{u} = \mathbf{f}$$

In order to get the equations of characteristic lines where PDE turns into ODEs, the eigenproblem

$$\mathbf{A}\partial_x\mathbf{u} = \frac{dX}{dt}\partial_x\mathbf{u} \,,$$

holds. This problem has non trivial solution if $\frac{dX}{dt}$ and $\partial_x\mathbf{u}$ are pairs of eigenvalues and (right) eigenvectors of the array $\mathbf{A}$.

**Diagonalization.**

$$\mathbf{A} = \mathbf{R}\Lambda\mathbf{L}$$

$$\mathbf{L}\left[\partial_t\mathbf{u} + \mathbf{R}\Lambda\mathbf{L}\partial_x\mathbf{u}\right] = \mathbf{L}\mathbf{f}$$

Since $\mathbf{L} = \mathbf{R}^{-1}$, and defining the **characteristic variables** by $d\mathbf{q} = \mathbf{L}d\mathbf{u}$ - in linear problems matrix $\mathbf{A}$ is constant, and so its spectral decomoposition, and thus $\mathbf{q} = \mathbf{L}\mathbf{u}$ - , it's possible to recast the original problem in diagonal form

$$\partial_t\mathbf{q} + \Lambda\partial_x\mathbf{q} = \mathbf{L}\mathbf{f}$$

$$\partial_t q_i + \Lambda_i\partial_x q_i = \sum_k L_{ik} f_k =: F_i \,.$$

Thus, on the $i^{th}$ family of characteristic lines, $\dfrac{dX}{dt} = \lambda_i$, $Q_i(t) = q_i(x(t), t)$ evolves as

$$\frac{dQ_i}{dt} = F_i \,.$$

If $F_i = [\mathbf{L}\mathbf{f}]_i = 0$, the characteristic variable $q_i$ is constant along the characteristic lines. Once the characteristic variables are determined, the conservative variables are evalauted as $\mathbf{u}(x, t) = \mathbf{R}\mathbf{q}(x, t)$.

## Domain of influence and domain of dependence

## Riemann problem

A Riemann problem is defined as the evolution of the initial state

$$\mathbf{u}(x, t_0) = \begin{cases} \mathbf{u}_a \,, & x < x_0 \\ \mathbf{u}_b \,, & x > x_0 \end{cases}$$

This problem is quite useful in quite a wide range of numerical methods for hyperbolic problems - Godunov schemes in Finite Volume Methods -, to evaluate the **boundary state** to be used numerical flux.

For linear problems, the matrix $\mathbf{A}$ is constant ad so it is its spectral decomposition, $\mathbf{A} = \mathbf{R}\Lambda\mathbf{L}$, and the solution of a Riemann problem of an homogeneous linear hyperbolic system can be easily determined analytically with the method of characteristics,

Let's change the origin of space and time, so that the initial state is in $t = 0$, and the jump in the initial condiiton in $x = 0$. Each charactersitic variable $q_k(x, t)$ is constant on its family of characteristic lines, $x = X_k(t) = x_{0,k} + \lambda_k t$.

$$q_k(x, t) = q_k(x_{0,k} + \lambda_k t, t) = q_k(x_{0,k}, 0) = q_k(x - \lambda_k t, 0) = L_{ki}u_j(x - \lambda_k t, 0) \,.$$

Thus, the solution in conservative variables $\mathbf{u}(x,t)$ in $x$ at time $t$ reads

$$\mathbf{u}(x,t) = \mathbf{R}\mathbf{q}(x,t)$$
$$u_i(x,t) = R_{ik}q_k(x,t) = R_{ik}q_k(x - \lambda_k t, 0) = R_{ik}L_{kj}u_j(x - \lambda_k t, 0)$$

In a Riemann problem for a $N$-dimensional linear system the solution shows $N+1$ homogeneous regions (at most, in general the same number as the number of the non-coincident eigenvalues $+1$), delimited by the characteristic lines with origin in the discontinuity. Sorting the eigenvalues in increasing order

$$\lambda_1 > \lambda_2 > \cdots > \lambda_N \ ,$$

and defining the homogeneous regions

$$S_0 : \frac{x}{t} \in (-\infty, \lambda_1)$$
$$S_1 : \frac{x}{t} \in (\lambda_1, \lambda_2)$$
$$...$$
$$S_i : \frac{x}{t} \in (\lambda_i, \lambda_{i+1})$$
$$...$$
$$S_{N-1} : \frac{x}{t} \in (\lambda_{N-1}, \lambda_N)$$
$$S_N : \frac{x}{t} \in (\lambda_N, +\infty)$$

the solution is in the $S_i$ region is

$$u_i(x,t) = \sum_{\lambda_k > \frac{x}{t}} R_{ik}q_{a,k} + \sum_{\lambda_k < \frac{x}{t}} R_{ik}q_{b,k}$$

---

**Example 28.3.1 (Linear(ized) P-system)**

The linear(ized) P-system around a uniform reference state $\overline{\rho}$, $\overline{u}$ in convective form reads

$$\partial_t \begin{bmatrix} \rho \\ u \end{bmatrix} + \begin{bmatrix} \overline{u} & \overline{\rho} \\ \frac{a^2}{\overline{\rho}} & \overline{u} \end{bmatrix} \partial_x \begin{bmatrix} \rho \\ u \end{bmatrix} = \mathbf{0} \ .$$

**Spectral decomposition.**

$$0 = |-\lambda\mathbf{I} + \mathbf{A}| = (\overline{u} - \lambda)^2 - a^2$$

$$\lambda_{12} = \overline{u} \mp a \qquad , \qquad \mathbf{r}_{12} = \begin{bmatrix} \overline{\rho} \\ \mp a \end{bmatrix}$$

$$\mathbf{R} = \begin{bmatrix} \overline{\rho} & \overline{\rho} \\ -a & a \end{bmatrix}$$

$$\mathbf{L} = \mathbf{R}^{-1} = \frac{1}{2\overline{\rho}a} \begin{bmatrix} a & -\overline{\rho} \\ a & \overline{\rho} \end{bmatrix}$$

**Reference state.**

$$|u| : \begin{cases} = 0 & \text{at rest} \\ < a & \text{subsonic flow} \\ > a & \text{supersonic flow to the left/right} \end{cases}$$

Subsonic: the two families of characteristic lines have opposite direction; supersonic: the two families of characteristic lines have the same direction.

**Example 28.3.2 (Linearized shallow water equations)**

**Example 28.3.3 (Linearized Euler equations (acoustics))**

**Example 28.3.4 (Wave equation)**

A wave equation arises in many different fields of science. As an example, 1-dimensional wave equation descrives the axial dynamics of a truss

$$m\partial_{tt}u - EA\partial_{xx}u = f \,,$$

that can be recast in the general expression of wave equation

$$\partial_{tt}u - c^2\partial_{xx}u = F$$

The $2^{nd}$ order differential operator appearing in 1-dimensional wave equation can be factored as the "product" of 2 $1^{st}$ order differentail operators,

$$\left(\partial_{tt} - c^2\partial_{xx}\right)u = \left(\partial_t - c\partial_x\right)\left(\partial_t + c\partial_x\right)u$$

and thus a wave equation can be written as

$$\begin{cases} \partial_t u + c\partial_x u - w = 0 \\ \partial_t w - c\partial_x w = F \end{cases}$$

In the regime of small displacement, the velocity field is the partial time derivative of the displacement field, $v = \partial_t u$, and the axial force reads $N = EA\partial_x u$. Exploiting Schwartz's theorem about mixed partial derivatives to write $\partial_t N = EA\partial_x v$, it's possible to write the wave function as the following system of hyperbolic equations in the physical unknowns $v, N$

$$\begin{cases} \partial_t N - EA\partial_x v = 0 \\ \partial_t v - \frac{1}{m}\partial_x N = f \end{cases}$$

**P-system and wave equation - reference state at rest, $\overline{u} = 0$.**

$$\begin{cases} \partial_t \rho + \overline{\rho}\partial_x u = 0 \\ \partial_t u + \frac{a^2}{\overline{\rho}}\partial_x \rho = 0 \end{cases}$$

Taking time partial derivative of the first and space partial derivative of the second equation times $\overline{\rho}$, and evaluating their difference, a wave equation for $rho$ appears

$$\partial_{tt}\rho - a^2\partial_{xx}\rho = 0 \,.$$

Analogously, taking space derivative of the first and time derivative of the second, a wave equation for the velocity field appears

$$\partial_{tt}u - a^2\partial_{xx}u = 0 \,.$$

## 28.4 System non-linear

### 28.4.1 1-dimensional space

$$\mathbf{u}(x,t)$$

$$\partial_t \mathbf{u} + \partial_x \mathbf{F}(\mathbf{u}) = \mathbf{f} \qquad \text{(conservative form)}$$
$$\partial_t \mathbf{u} + \partial_\mathbf{u} \mathbf{F}(\mathbf{u}) \partial_x \mathbf{u} = \mathbf{f} \quad \text{(convective form)}$$

## 28.5 n-dimensional space

$$\mathbf{u}(\vec{r},t)$$

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{F}(\mathbf{u}) = \mathbf{f} \qquad \text{(conservative form)}$$
$$\partial_t \mathbf{u} + \nabla \mathbf{u} \cdot \partial_\mathbf{u} \mathbf{F}(\mathbf{u}) = \mathbf{f} \quad \text{(convective form)}$$

Different descritpions of integral problem,

$$\frac{d}{dt} \int_V \mathbf{u} + \oint_{\partial V} \hat{n} \cdot \mathbf{F}(\mathbf{u}) = \int_V \mathbf{f} \qquad \text{(Eulerian)}$$

$$\frac{d}{dt} \int_{V_t} \mathbf{u} - \oint_{\partial V_t} \mathbf{u} \vec{u} \cdot \hat{n} + \oint_{\partial V_t} \hat{n} \cdot \mathbf{F}(\mathbf{u}) = \int_{V_t} \mathbf{f} \qquad \text{(Lagrangian)}$$

$$\frac{d}{dt} \int_{v_t} \mathbf{u} - \oint_{\partial v_t} \mathbf{u} \vec{u}_b \cdot \hat{n} + \oint_{\partial v_t} \hat{n} \cdot \mathbf{F}(\mathbf{u}) = \int_{v_t} \mathbf{f} \qquad \text{(arbitrary)}$$

**todo** *in coordinates*

$$f_i = \partial_t u_i + \partial_{x_k} F_{ki}(u_l) =$$
$$= \partial_t u_i + \partial_{x_k} u_m \partial_{u_m} F_{ki}(u_l) =$$

---

**Example 28.5.1 (P-system in 1-dimensional domain)**

$$\begin{cases} \partial_t \rho + u \partial_x \rho + \rho \partial_x u = 0 \\ \rho \partial_t u + \rho u \partial_x u + \partial_x P = 0 \end{cases}$$

with $\partial_x P = a^2 \partial_x \rho$,

**Convective form**

$$\partial_t \begin{bmatrix} \rho \\ u \end{bmatrix} + \begin{bmatrix} u & \rho \\ \frac{a^2}{\rho} & u \end{bmatrix} \partial_x \begin{bmatrix} \rho \\ u \end{bmatrix} = \underline{0} \,.$$

**Conservative form**

$$\partial_t \begin{bmatrix} \rho \\ \rho u \end{bmatrix} + \partial_x \begin{bmatrix} \rho u \\ \rho u^2 + \rho a^2 \end{bmatrix} = \underline{0} \,.$$

**Spectral decomposition** of $\mathbf{A}(\mathbf{u})$ gives

$$0 = \left| \left| \begin{bmatrix} u - s & \rho \\ \frac{a^2}{\rho} & u - s \end{bmatrix} \right| \right| = (u - s)^2 - a^2$$

$$s_{1,2} = u \mp a$$

$$\mathbf{R} = \begin{bmatrix} \rho & \rho \\ a & -a \end{bmatrix}$$

$$\mathbf{L} = \frac{1}{2\rho a} \begin{bmatrix} a & \rho \\ a & -\rho \end{bmatrix}$$

**Example 28.5.2 (Euler equations in 1-dimensional domain)**

**Conservative form**

$$\partial_t \begin{bmatrix} \rho \\ \rho u \\ \rho e^t \end{bmatrix} + \partial_x \begin{bmatrix} \rho u \\ \rho u^2 + P \\ \rho u h^t \end{bmatrix} = \underline{0} \,,$$

with $h^t = e^t + \frac{P}{\rho}$ and $e^t = e + \frac{u^2}{2}$, and the pressure field can be written as a function of the other thermodynamic variables. As an example, using conservative variables $(\rho, m, E^t) = (\rho, \rho u, \rho e^t) = \left( \rho, \rho u, \rho \left( e + \frac{u^2}{2} \right) \right)$

$$P(\rho, e) = P \left( \rho, \frac{E^t}{\rho} - \frac{m^2}{2\rho^2} \right) = \Pi \left( \rho, m, E^t \right)$$

so that

$$\partial_\rho \Pi = \partial_\rho P\big|_e + \partial_e P\big|_\rho \left( -\frac{E^t}{\rho^2} + \frac{m^2}{\rho^3} \right) =$$

$$= \partial_\rho P\big|_e + \partial_e P\big|_\rho \left( -\frac{e^t}{\rho} + \frac{u^2}{\rho} \right)$$

$$= c^2 - \frac{P}{\rho^2} \partial_e P\big|_\rho + \partial_e P\big|_\rho \left( -\frac{e^t}{\rho} + \frac{u^2}{\rho} \right)$$

$$= c^2 + \partial_e P\big|_\rho \left( -\frac{h^t}{\rho} + \frac{u^2}{\rho} \right)$$

$$\partial_m \Pi = \partial_e P\big|_\rho \left( -\frac{m}{\rho^2} \right)$$

$$\partial_{E^t} \Pi = \partial_e P\big|_\rho \left( \frac{1}{\rho} \right)$$

The speed of sound reads

$$c^2 = \partial_\rho P\big|_s =$$

$$= \partial_\rho P\big|_e + \partial_e P\big|_\rho \, \partial_\rho e\big|_s =$$

$$= \partial_\rho P\big|_e + \frac{P}{\rho^2} \partial_e P\big|_\rho \,,$$

**Conservative form in conservative variables.**

$$\partial_t \begin{bmatrix} \rho \\ m \\ E^t \end{bmatrix} + \partial_x \begin{bmatrix} m \\ \frac{m^2}{\rho} + \Pi(\rho, m, E^t) \\ \frac{m}{\rho} \left( E^t + \Pi(\rho, m, E^t) \right) \end{bmatrix} = \underline{0} \,,$$

**Convective form in conservative variables.**

$$\partial_t \begin{bmatrix} \rho \\ m \\ E^t \end{bmatrix} + \partial_x \begin{bmatrix} 0 & 1 & 0 \\ -\frac{m^2}{\rho^2} + \partial_\rho \Pi & \frac{2m}{\rho} + \partial_m \Pi & \partial_{E^t} \Pi \\ -\frac{m}{\rho^2}(E^t + \Pi) + \frac{m}{\rho} \partial_\rho \Pi & \frac{1}{\rho}(E^t + \Pi) + \frac{m}{\rho} \partial_m \Pi & \frac{m}{\rho} \left( 1 + \partial_{E^t} \Pi \right) \end{bmatrix} \partial_x \begin{bmatrix} \rho \\ m \\ E^t \end{bmatrix} = \underline{0} \,,$$

**Spectral decomposition** of $\mathbf{A}(\mathbf{u})$

$$0 = \left| \left[ \begin{array}{ccc} -s & 1 & 0 \\ -u^2 + \partial_\rho \Pi & 2u + \partial_m \Pi - s & \partial_{E^t} \Pi \\ -u\left(e^t + \frac{P}{\rho}\right) + u\partial_\rho \Pi & e^t + \frac{P}{\rho} + u\partial_m \Pi & u\left(1 + \partial_{E^t}\Pi\right) - s \end{array} \right] \right| =$$

$$= -s\left[ (2u + \partial_m \Pi - s)(u(1 + \partial_{E^t}\Pi) - s) - \partial_{E^t}\Pi (h^t + u\partial_m \Pi) \right] +$$
$$- uh^t \partial_{E^t}\Pi + u\partial_\rho \Pi \partial_{E^t}\Pi +$$
$$+ (u^2 - \partial_\rho \Pi)(u(1 + \partial_{E^t}\Pi) - s) =$$
$$= -s^3 +$$
$$+ s^2 \left( 2u + \partial_m \Pi + u + u\partial_{E^t}\Pi \right) +$$
$$+ s \left( -2u^2 - 2u^2 \partial_{E^t}\Pi - u\partial_m \Pi - u\partial_m \Pi \partial_{E^t}\Pi + \partial_{E^t}\Pi h^t + u\partial_{E^t}\Pi \partial_m \Pi - u^2 + \partial_\rho \Pi \right) +$$
$$+ \left( -uh^t \partial_{E^t}\Pi + u\partial_\rho \Pi \partial_{E^t}\Pi + u^3 + u^3 \partial_{E^t}\Pi - u\partial_\rho \Pi - u\partial_\rho \Pi \partial_{E^t}\Pi \right) +$$
$$= -s^3 +$$
$$+ s^2 \left( 3u + \partial_m \Pi + u\partial_{E^t}\Pi \right) +$$
$$+ s \left( -3u^2 - 2u^2 \partial_{E^t}\Pi - u\partial_m \Pi + \partial_{E^t}\Pi h^t + \partial_\rho \Pi \right) +$$
$$+ \left( u^3 - uh^t \partial_{E^t}\Pi + u^3 \partial_{E^t}\Pi - u\partial_\rho \Pi \right) =$$
$$= -(s - u)^3 + (s - u)c^2 =$$
$$= (s - u)\left[ -(s - u)^2 + c^2 \right]$$

being

$$\partial_m \Pi + u\partial_{E^t}\Pi = \left( -\frac{u}{\rho} + \frac{u}{\rho} \right) \partial_\rho P\big|_e = 0$$

$$-2u^2 \partial_{E^t}\Pi - u\partial_m \Pi + \partial_{E^t}\Pi h^t + \partial_\rho \Pi = -\frac{u^2}{\rho}\partial_\rho P\big|_e + \frac{1}{\rho}\partial_e P\big|_\rho h^t + \partial_\rho P\big|_e + \partial_e P\big|_\rho \left( -\frac{e^t}{\rho} + \frac{u^2}{\rho} \right) =$$
$$= \partial_e P\big|_\rho \frac{P}{\rho^2} + \partial_\rho P\big|_e =$$
$$= c^2$$

$$-u\partial_\rho \Pi + u^3 \partial_{E^t}\Pi - uh^t \partial_{E^t}\Pi = u\left( -\partial_\rho P\big|_e - \partial_e P\big|_\rho \left( -\frac{e^t}{\rho} + \frac{u^2}{\rho} \right) + \frac{u^2}{\rho}\partial_e P\big|_\rho - \frac{h^t}{\rho}\partial_e P\big|_\rho \right)$$
$$= u\left( -\partial_\rho P\big|_e - \frac{P}{\rho^2}\partial_e P\big|_\rho \right) =$$
$$= -uc^2 .$$

Thus,

$$s_{1,3} = u \mp c \quad , \quad s_2 = u$$

$$\mathbf{r}_{1,3} = \left[ \begin{array}{c} 1 \\ u \mp c \\ ... \end{array} \right] \hat{\rho} \quad , \quad \mathbf{r}_2 = \left[ \begin{array}{c} ... \\ 0 \\ ... \end{array} \right] \hat{\rho}$$

being

$$\hat{E}^t \partial_{E^t}\Pi = \left[ u^2 - \partial_\rho \Pi + (-u \pm c - \partial_m \Pi)(u \mp c) \right] \hat{\rho} =$$
$$= \left[ u^2 - \partial_\rho \Pi - u^2 + \pm 2uc - c^2 - \partial_m \Pi(u \mp c) \right] \hat{\rho} =$$
$$= \left[ -c^2 - \partial_e P\big|_\rho \left( -\frac{h^t}{\rho} + \frac{u^2}{\rho} \right) \pm 2uc - c^2 + \frac{u}{\rho}\partial_e P\big|_\rho (u \mp c) \right] \hat{\rho} =$$
$$= \left[ -2c^2 + \partial_e P\big|_\rho \frac{h^t}{\rho} \pm 2uc \mp \frac{uc}{\rho}\partial_e P\big|_\rho \right] \hat{\rho} =$$

$$\mathbf{R} = \dots$$

$$\mathbf{L} = \dots$$

---

**Example 28.5.3 (Shallow water equation in $1$-dimensional domain)**

Let $b(x)\dots$, $h(x)$ the height of the free surface, $\eta(x) = h(x) - b(x)$ the depth.

Derivative of integrals with non-constant extremes

$$\partial_x \int_{z=0}^{\eta(x,t)} \rho u \, dz = \int_{z=0}^{\eta(x,t)} \partial_x(\rho u) \, dz + \rho u(x, \eta(x,t), t) \partial_x \eta(x,t) \ .$$

Continuity equation reads

$$\partial_t \rho + \partial_x(\rho u) + \partial_z(\rho w) = 0 \ ,$$

for fluids with constant and uniform density

$$0 = \int_{z=0}^{\eta(x,t)} \left( \partial_t \rho + \partial_x(\rho u) + \partial_z(\rho w) \right) \, dz =$$

$$= \partial_x \int_{z=0}^{\eta(x,t)} (\rho u) \, dz - \rho u(x, \eta, t) \partial_x \eta + \rho w(x, \eta(x,t), t) =$$

$$= \partial_x \int_{z=0}^{\eta(x,t)} (\rho u) \, dz + \rho \partial_t \eta = \qquad\qquad \simeq \partial_x (\rho \eta u) + \partial_t (\rho \eta) \ .$$

having linked the velocity to the material derivative of the position, whose vertical component reads

$$w(x, \eta(x,t), t) = \frac{D\eta}{Dt} = \partial_t \eta(x,t) + u(x, \eta(x,t), t) \partial_x \eta \ .$$

Assuming hydrostatic pressure distribution, $p = P_a + \rho g z$ at depth $z$ under the level of local free surface,

Momentum equation reads

$$0 = \partial_t(\rho u) + \partial_x(\rho u^2) + \partial_z(\rho u w) + \partial_x P \ .$$

and integration in $z$-direction **todo** Explicitly treat the $z$ term

$$0 = \partial_t(\rho \eta u) + \partial_x(\rho u^2 \eta) + \partial_x \int_{z=0}^{\eta(x)} (P_a + \rho g z) \, dz =$$

$$= \partial_t(\rho \eta u) + \partial_x \left( \rho u^2 \eta + \frac{1}{2} \rho g \eta^2 \right) \ .$$

**Conservative form of the equations.**

$$\begin{cases} \partial_t(\eta) + \partial_x m = 0 \\ \partial_t m + \partial_x \left( \frac{m^2}{\eta} + \frac{g\eta^2}{2} \right) = 0 \end{cases}$$

**Convective form of the equations.**

$$\partial_t \begin{bmatrix} \eta \\ m \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ -\frac{m^2}{\eta^2} + g\eta & 2\frac{m}{\eta} \end{bmatrix} \partial_x \begin{bmatrix} \eta \\ m \end{bmatrix} = \underline{0}$$

---

**Spectrum of matrix A(u).**

$$0 = |\mathbf{A}(\mathbf{u}) - s^2\mathbf{I}| = -s\,(2u - s) + u^2 - g\eta = (s - u)^2 - g\eta \; .$$

---

**Example 28.5.4 (P-system in n-dimensional domain)**

- Conservative variables: $(\rho, \vec{m})$

- Physical variables: e.g. $(\rho, \vec{u})$, $(P, \vec{u})$,…

$$\begin{cases} \partial_t \rho + \nabla \cdot \vec{m} = 0 \\ \partial_t \vec{m} + \nabla \cdot \left[ \frac{\vec{m} \otimes \vec{m}}{\rho} + \rho a^2 \mathbb{I} \right] = 0 \end{cases}$$

---

**Example 28.5.5 (Euler system in n-dimensional domain)**

- Conservative variables: $(\rho, \vec{m}, E^t)$

- Physical variables: e.g. $(\rho, \vec{u}, e)$,…

$$\begin{cases} \partial_t \rho + \nabla \cdot \vec{m} = 0 \\ \partial_t \vec{m} + \nabla \cdot \left[ \frac{\vec{m} \otimes \vec{m}}{\rho} + \Pi\,\mathbb{I} \right] = \vec{0} \\ \partial_t E^t + \nabla \cdot \left[ \frac{\vec{m}(E^t + \Pi)}{\rho} \right] = 0 \end{cases}$$

where $\Pi$ represents the pressure field as a function of the conservative varaibles,

$$\Pi\,(\rho, \vec{m}, E^t) = P\,(\rho, e) = P\left( \rho, \frac{E^t}{\rho} - \frac{|\vec{m}|^2}{\rho^3} \right) \; ,$$

and $P$ the pressure field expressed by the **equation of state of the fluid** as a function of density and internal energy per unit mass as the pair of independent variables determining the thermodynamic state.

---

**Example 28.5.6 (Shallow water equations in 2-dimensional domain)**

$$\begin{cases} \partial_t (\rho\eta) + \nabla \cdot (\rho\eta\vec{u}) = 0 \\ \partial_t (\rho\eta\vec{u}) + \nabla \cdot \left( \rho\eta\vec{u}\vec{u} + \frac{1}{2}\rho g\eta^2 \mathbb{I} \right) = 0 \end{cases}$$

# 28.6  Method of characteristics

## 28.6.1  Scalar multi-dimensional problem

$$\mathbf{a}^T\,(u(\mathbf{x}), \mathbf{x})\,\nabla u(\mathbf{x}) = c(u(\mathbf{x}), \mathbf{x}) \quad , \quad \mathbf{x} \in D$$

$$a_k(u(\mathbf{x}), \mathbf{x})\frac{\partial}{\partial x_k}u(\mathbf{x}) = c(u(\mathbf{x}), \mathbf{x})$$

The method of characteristics aims at transforming the PDE over the domain $D$, in a set of ODEs over lines - the characteristic lines. Both the coordinates of these lines and the unknown function on these lines are written in parametric form as

$$\mathbf{x} = \mathbf{X}(s)$$
$$u(\mathbf{X}(s)) = U(s) \,.$$

Now, taking the ordinary derivative of $U'(s)$,

$$U'(s) = \frac{dX_k}{ds}(s) \frac{\partial u}{\partial x_k}(\mathbf{X}(s)) \,.$$

If the characteristic lines are defined by the differential equation $X_k'(s) = a_k(U(s), X(s))$, the set of ODEs become

$$U'(s) = c\left(U(s), \mathbf{X}(s)\right)$$
$$\mathbf{X}'(s) = \mathbf{a}\left(U(s), \mathbf{X}(s)\right) \,.$$

## 28.6.2 Vector multi-dimensional problem

**todo** *Characteristic method fails…characteristic - eikonal - equation. Uncomment*

Let

$$\mathbf{u}(t, \mathbf{r}) : \mathbb{R} \times \mathbb{R}^d \to \mathbb{R}^n$$
$$\mathbf{F}(\mathbf{u}) : \mathbb{R}^n \to \mathbb{R}^d \times \mathbb{R}^n$$
$$\mathbf{s}(t, \mathbf{r}) : \mathbb{R} \times \mathbb{R}^d \to \mathbb{R}^n$$

the conservative form of a hyperbolic problem reads

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{F}(\mathbf{u}) = \mathbf{s}$$
$$\partial_t u_i + \partial_j F_{ji} = s_i \,.$$

Expanding the divergence of the flux, the convective (quasi-linear? It makes little to no sense to me…) form follows

$$s_i = \partial_t u_i + \sum_{j=1}^{d} \partial_j F_{ji} =$$

$$= \partial_t u_i + \sum_{j=1}^{d} \sum_{k=1}^{n} \partial_j u_k \partial_{u_k} F_{ji} =$$

$$= \partial_t u_i + \sum_{j=1}^{d} \sum_{k=1}^{n} A_{ik}^{j}(\mathbf{u}) \, \partial_j u_k \,,$$

for every component $i = 1 : n$, or

$$\mathbf{s} = \partial_t \mathbf{u} + \sum_{j=1}^{d} \mathbf{A}^j \partial_j \mathbf{u} \,.$$

This problem can be made a little more general with by defining $\mathbf{x} = (t, \mathbf{r}) = (x_0, \mathbf{r})$, as

$$\sum_{j=0}^{d} \mathbf{A}^j \, \partial_j \mathbf{u} = \mathbf{s} \,.$$

The former expression of the hyperbolic problem immediately follows if $\mathbf{A}_0 = \mathbf{I}_n$, $A_{ik}^0 = \delta_{ik}$, and

$$\widetilde{\nabla} \cdot \widetilde{\mathbf{F}} = \partial_0 \mathbf{F}_0 + \sum_{j=1}^{d} \partial_j \mathbf{F}_j = \widetilde{\nabla} \cdot [\, \mathbf{u} \mid \mathbf{F} \,] \,.$$

as, for $i = 1 : n$,

$$\delta_{ik} = \partial_{u_k} F_{0i} \quad \rightarrow \quad F_{0i} = u_k \delta_{ik} = u_i \,.$$

## Taylor expansion of a solution

Starting from the solution on a manifold determined by the equation $S : f(\mathbf{x}) = 0$, $\mathbf{x}_0 \in S$, if the solution is differentiable, the solution in a point $\mathbf{x} = \mathbf{x}_0 + \Delta\mathbf{x}$ reads

$$\mathbf{u}(\mathbf{x}) \sim \mathbf{u}(\mathbf{x}_0) + \Delta\mathbf{x} \cdot \nabla\mathbf{u}(\mathbf{x}_0)$$

$$u_i(\mathbf{x}) \sim u_i(\mathbf{x}_0) + \Delta_j \frac{\partial u_i}{\partial x_j} \ .$$

Now, with a change of coordinates from $\mathbf{x}$ to $\xi = (n, \mathbf{t})$, i.e. to local normal and tangential direction,

$$\nabla\mathbf{u} = \hat{\mathbf{x}}_i \frac{\partial}{\partial x_i}\mathbf{u} = \hat{\mathbf{x}}_i \frac{\partial \mathbf{u}}{\partial x_i} = \hat{\mathbf{x}}_i \frac{\partial \mathbf{u}}{\partial \xi_k}\frac{\partial \xi_k}{\partial x_i} = \hat{\xi}_k \frac{\partial \mathbf{u}}{\partial \xi_k}$$

Inserting in the hyperbolic equation

$$s_i = \sum_{j=0}^{d} A_{ik}^j \frac{\partial u_k}{\partial x_j} =$$

$$= \sum_{j=0}^{d} A_{ik}^j \sum_{\ell=0}^{d} \frac{\partial u_k}{\partial \xi_\ell}\frac{\partial \xi_\ell}{\partial x_j} =$$

$$= \sum_{\ell=0}^{d}\sum_{j=0}^{d} A_{ik}^j \frac{\partial \xi_\ell}{\partial x_j}\frac{\partial u_k}{\partial \xi_\ell} =$$

$$= \sum_{\ell=0}^{d}\sum_{j=0}^{d} A_{ik}^j \xi_j^\ell \frac{\partial u_k}{\partial \xi_\ell} \ ,$$

and separating the normal $\ell = 0$ from the tangential $\ell = 1 : d$ coordinates,

$$s_i = \sum_{j=0}^{d} A_{ik}^j \xi_j^0 \frac{\partial u_k}{\partial \xi_0} + \sum_{\ell=1}^{d}\sum_{j=0}^{d} A_{ik}^j \xi_j^t \frac{\partial u_k}{\partial \xi_t} =$$

$$= \sum_{j=0}^{d} A_{ik}^j n_j \frac{\partial u_k}{\partial n} + \sum_{\ell=1}^{d}\sum_{j=0}^{d} A_{ik}^j t_j^\ell \frac{\partial u_k}{\partial t_\ell} =$$

$$= \sum_{j=0}^{d} n_j \mathbf{A}^j \partial_n \mathbf{u} + \sum_{\ell=1}^{d}\sum_{j=0}^{d} t_j^\ell \mathbf{A}^j \partial_{t_\ell} \mathbf{u} \ .$$

On a smooth surface where the solution is known, all the tangential derivatives of the solution are known as well. In order to evaluate the normal derivative $\partial_n \mathbf{u}$ from the PDE, the (formal) inversion of the matrix

$$\mathbf{A}_{\hat{\mathbf{n}}} := \sum_{j=0}^{d} n_j \, \mathbf{A}^j \ ,$$

must be invertible, to formally get

$$\partial_n \mathbf{u} = \mathbf{A}_{\hat{\mathbf{n}}}^{-1}\left(\mathbf{s} - \sum_{\ell=1}^{d} \mathbf{A}_{\hat{\mathbf{t}}_\ell} \partial_{t_\ell} \mathbf{u}\right) \ .$$

This expression of the normal derivative of the solution - whenever it exists - can be used to find the approximation of the solution in normal direction w.r.t. the surface $S$, i.e. with $\Delta\mathbf{x} = \Delta\ell\hat{\mathbf{n}}$ as

$$\mathbf{u}(\mathbf{x}) \sim \mathbf{u}(\mathbf{x}_0) + \Delta\ell\,\hat{\mathbf{n}} \cdot \nabla\mathbf{u}(\mathbf{x}_0) =$$

$$= \mathbf{u}(\mathbf{x}_0) + \Delta\ell\,\partial_n \mathbf{u}(\mathbf{x}_0) \ .$$

**If** the matrix $\mathbf{A}_{\hat{\mathbf{n}}}$ is diagonalizable, it's invertible if it has no zero eigenvalue. Let $\hat{\mathbf{n}}_i$ be a unit vector so that the matrix $\mathbf{A}_{\hat{\mathbf{n}}}$ has an eigenvalue equal to zero $s_0 = 0$, with right and left eigenvectors $\mathbf{r}_0$, $\mathbf{l}_0$. Recalling the expression of the PDE in normal and tangential components

$$\mathbf{s} = \mathbf{A}_{\hat{\mathbf{n}}}\partial_n\mathbf{u} + \sum_{\ell=1}^{d} \mathbf{A}_{\hat{\mathbf{d}}_\ell}\partial_{t_\ell}\mathbf{u} \ ,$$

left-multiplying by $\mathbf{l}$ gives

$$\mathbf{l}_0^T\mathbf{s} = \mathbf{l}_0^T\sum_{\ell=1}^{d} \mathbf{A}_{\hat{\mathbf{d}}_\ell}\partial_{t_\ell}\mathbf{u} = \qquad \cdots \qquad = \mathbf{l}_0^T\sum_{k=0}^{d} \mathbf{A}^k\partial_k\mathbf{u} \ .$$

**todo**

- *compatibility conditions*

- *eikonal equation*

- *explicitly treating steady vs. unsteady problems*

### Examples

### Isothermal compressible flow

- Conservative variables: $(\rho, \vec{m})$

- Physical variables: e.g. $(\rho, \vec{u})$, $(P, \vec{u})$,…

Conservative form reads

$$\begin{cases} \partial_t\rho + \nabla \cdot \mathbf{m} = 0 \\ \partial_t\mathbf{m} + \nabla \cdot \left[\frac{\mathbf{m}\otimes\mathbf{m}}{\rho} + \rho a^2\mathbf{I}\right] = 0 \end{cases}$$

Convective form in a 2-dimensional domain reads

$$\partial_t \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} + \begin{bmatrix} \cdot & 1 & \cdot \\ a^2 - \frac{m_x m_x}{\rho^2} & \frac{2m_x}{\rho} & \cdot \\ -\frac{m_x m_y}{\rho^2} & \frac{m_y}{\rho} & \frac{m_x}{\rho} \end{bmatrix} \partial_x \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} + \begin{bmatrix} \cdot & \cdot & 1 \\ -\frac{m_x m_y}{\rho^2} & \frac{m_y}{\rho} & \frac{m_x}{\rho} \\ a^2 - \frac{m_y m_y}{\rho^2} & \cdot & \frac{2m_y}{\rho} \end{bmatrix} \partial_y \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} = \mathbf{0}$$

### Some algebra

$$0 = \rho_{/t} + m_{x/x} + m_{y/y}$$
$$0 = m_{i/t} + \partial_j\left(\frac{m_j m_i}{\rho} + \rho a^2\delta_{ij}\right) \ .$$

Thus,

$$\begin{aligned} \mathbf{A}_{\hat{\mathbf{n}}} &= \begin{bmatrix} \cdot & n_x & n_y \\ a^2 n_x - u_n u_x & u_x n_x + u_n & u_x n_y \\ a^2 n_y - u_n u_y & u_y n_x & u_n + u_y n_y \end{bmatrix} = \\ &= \begin{bmatrix} 0 & \mathbf{n}^T \\ a^2\mathbf{n} - (\mathbf{n}^T\mathbf{u})\mathbf{u} & \mathbf{n}^T\mathbf{u}\mathbf{I} + \mathbf{u}\mathbf{n}^T \end{bmatrix} = \\ &= \begin{bmatrix} 0 & \hat{\mathbf{n}}^T \\ a^2\mathbf{n} - (\mathbf{n} \cdot \mathbf{u})\mathbf{u} & (\mathbf{n} \cdot \mathbf{u})\mathbb{I} + \mathbf{u} \otimes \mathbf{n} \end{bmatrix} \ . \end{aligned}$$

## Eigenvalues. Details

The eigenvalue decomposition of the matrix $\mathbf{A}_{\hat{\mathbf{n}}}$ follows from

$$0 = |\mathbf{A}_{\hat{\mathbf{n}}} - s\mathbf{I}| = \left|\begin{bmatrix} -s & \mathbf{n}^T \\ a^2\mathbf{n} - u_n\mathbf{u} & (u_n - s)\mathbf{I} + \mathbf{u}\mathbf{n}^T \end{bmatrix}\right| =$$

$$= -s(u_n + u_x n_x - s)(u_n + u_y n_y - s) + n_x u_x n_y(a^2 n_y - u_n u_y) + n_y u_y n_x(a^2 n_x - u_n u_x) +$$
$$- (a^2 n_y - u_n u_y)(u_x n_x + u_n - s)n_y - (a^2 n_x - u_n u_x)(u_y n_y + u_n - s)n_x - s u_x u_y n_x n_y =$$
$$= s^3(-1) +$$
$$\quad + s^2(u_n + u_y n_y + u_n + u_x n_x) +$$
$$\quad + s(-(u_n + u_x n_x)(u_n + u_y n_y) + n_y(a^2 n_y - u_n u_y) + n_x(a^2 n_x - u_n u_x) - u_x n_x u_y n_y) +$$
$$\quad + 1 \cdot (n_x u_x n_y(a^2 n_y - u_n u_y) + n_y u_y n_x(a^2 n_x - u_n u_x) - n_y(a^2 n_y - u_n u_y)(u_x n_x + u_n) - n_x(a^2 n_x - u_n u_x)(u_y n_y + u_n)) =$$
$$= -s^3 + 3u_n s^2 + s(-u_n^2 - u_n(u_x n_x + u_y n_y) - u_x n_x u_y n_y - u_n u_y n_y - u_n u_x n_x + a^2 - u_x n_x u_y n_y) +$$
$$\quad + (a^2\left(n_y^2 u_x n_x + n_x^2 n_y u_y - n_y^2 u_x n_x - n_y^2 u_n - n_x^2 u_y n_y - u_n n_x^2\right) - u_n u_x n_x u_y n_y - u_n u_x n_x u_y n_y + u_n u_x n_x u_y n_y + u_n^2 u_y n_y +$$
$$= -s^3 + 3u_n s^2 + s(a^2 - 3u_n^2) - u_n(a^2 - u_n^2) =$$
$$= -(s - u_n)^3 + (s - u_n)a^2 =$$
$$= -(s - u_n)((s - u_n)^2 - a^2) .$$

Thus the eigenvalues of the matrix $\mathbf{A}_{\hat{\mathbf{n}}}$ are

$$s_{1,3} = u_n \mp a \quad , \quad s_2 = u_n ,$$

and the right and left eigenvectors follow

$$\mathbf{R} =$$
$$\mathbf{L} =$$

## Right and left eigenvectors. Details

For $s_{1,3} = u_n \mp a$,

$$\mathbf{0} = \begin{bmatrix} -u_n \pm a & \mathbf{n}^T \\ a^2\mathbf{n} - u_n\mathbf{u} & \pm a\mathbf{I} + \mathbf{u}\mathbf{n}^T \end{bmatrix}\begin{bmatrix} \hat{\rho} \\ \hat{\mathbf{m}} \end{bmatrix}$$

From the first equation it follows $\hat{m}_n = (u_n \mp a)\hat{\rho}$, and thus the momentum equation gives

$$\mathbf{0} = \hat{\rho}(a^2\mathbf{n} - u_n\mathbf{u}) \pm a\hat{\mathbf{m}} + \mathbf{u}\hat{m}_n =$$
$$= \hat{\rho}(a^2\mathbf{n} - u_n\mathbf{u}) \pm a\hat{\mathbf{m}} + \mathbf{u}(u_n \mp a)\hat{\rho} =$$
$$= \hat{\rho}(a^2\mathbf{n} \mp a\mathbf{u}) \pm a\hat{\mathbf{m}} ,$$

and thus

$$\hat{\mathbf{m}}_{1,3} = \hat{\rho}_{1,3}(\mathbf{u} \mp a\mathbf{n}) .$$

For $s_2 = u_n$,

$$\mathbf{0} = \begin{bmatrix} -u_n & \mathbf{n}^T \\ a^2\mathbf{n} - u_n\mathbf{u} & \mathbf{u}\mathbf{n}^T \end{bmatrix}\begin{bmatrix} \hat{\rho} \\ \hat{\mathbf{m}} \end{bmatrix}$$

From the first equation it follows $\hat{m}_n = u_n\hat{\rho}$, and thus the momentum equation gives

$$\mathbf{0} = \hat{\rho}(a^2\mathbf{n} - u_n\mathbf{u}) + \mathbf{u}\hat{m}_n =$$
$$= \hat{\rho}(a^2\mathbf{n} - u_n\mathbf{u}) + \mathbf{u}u_n\hat{\rho} =$$
$$= \hat{\rho}a^2\mathbf{n} ,$$

and thus $\hat{\rho}_2 = 0$. The components of the momentum readily follows from mass equation that is equivalent to the condition

$$\mathbf{n}^T \hat{\mathbf{m}}_2 = 0 \,,$$

i.e. as an example $\hat{\mathbf{m}}_2 = \begin{bmatrix} n_y \\ -n_x \end{bmatrix}$.

Assuming positive density, it's possible to write all the eigenvectors with the proper physical dimensions, and collect them in the matrix of right eigenvectors,

$$\mathbf{R} = \begin{bmatrix} \rho & 0 & \rho \\ \rho(u - an_x) & \rho an_y & \rho(u + an_x) \\ \rho(v - an_y) & -\rho an_x & \rho(v + an_y) \end{bmatrix} \,.$$

The determinant reads

$$\frac{1}{\rho^3 a^2} |\mathbf{R}| = n_y(v + n_y) - (u - n_x)n_x - n_y(v - n_y) + n_x(u + n_x) = 2 \,.$$

The inverse matrix reads

$$\mathbf{L} = \frac{1}{|\mathbf{R}|} \rho^2 \begin{bmatrix} a(a + u_n) & -an_x & -an_y \\ 2a(n_x v - n_y u) & 2an_y & -2an_x \\ a(a - u_n) & an_x & an_y \end{bmatrix} =$$

$$= \frac{1}{2\rho a^2} \begin{bmatrix} a(a + u_n) & -an_x & -an_y \\ 2a(n_x v - n_y u) & 2an_y & -2an_x \\ a(a - u_n) & an_x & an_y \end{bmatrix} \,.$$

as

$$L_{11} \propto an_y(v + an_y) + an_x(u + an_x) = au_n + a^2$$
$$-L_{21} \propto (u - an_x)(v + an_y) - (v - an_y)(u + an_x) = 2aun_y - 2an_x v$$
$$L_{31} \propto -an_x(u - an_x) - an_y(v - an_y) = -au_n + a^2$$
$$-L_{12} \propto an_x$$
$$L_{22} \propto (v + an_y) - (v - an_y) = 2an_y$$
$$-L_{32} \propto -an_x$$
$$L_{13} \propto -an_y$$
$$-L_{23} \propto (u + an_x) - (u - an_x) = 2an_x$$
$$L_{33} \propto an_y$$

Some check

$$\mathbf{LR} = \frac{1}{2\rho a^2} \begin{bmatrix} a(a + u_n) & -a\mathbf{n}^T \\ 2a(n_x v - n_y u) & 2a\mathbf{t}^T \\ a(a - u_n) & a\mathbf{n}^T \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ \mathbf{u} - a\mathbf{n} & a\mathbf{t} & \mathbf{u} + a\mathbf{n} \end{bmatrix} \rho =$$

$$= \frac{1}{2a^2} \begin{bmatrix} a^2 + au_n - au_n + a^2 & -a\mathbf{n}^T \mathbf{t} & a^2 + au_n - au_n - a^2 \\ 2a(n_x v - n_y u) + 2a(n_y u - n_x v) & 2a^2 & 2a(n_x v - n_y u) + 2a(n_y u - n_x v) \\ a^2 - au_n + au_n - a^2 & a\mathbf{n}^T \mathbf{t} & a^2 - au_n + au_n + a^2 \end{bmatrix} = \mathbf{I}_3 \,.$$

**Normal vectors of the characteristic surfaces.** For **locally supersonic flows**, $|\mathbf{u}| > a$, there are three directions, i.e. three unit vectors $\hat{\mathbf{n}}_i$ that make the eigenvalues equal to zero,

$$\hat{\mathbf{n}}_{1,3} \cdot \mathbf{u} = \pm a \quad , \quad \hat{\mathbf{n}}_2 \cdot \mathbf{u} = 0 \,,$$

or equivalently

$$\hat{\mathbf{n}}_{1,3} \cdot \left( \mathbf{u} \mp a\hat{\mathbf{n}}_{1,3} \right) = 0 \quad , \quad \hat{\mathbf{n}}_2 \cdot \mathbf{u} = 0 \,.$$

Let $\hat{\mathbf{u}}$ the unit vector along the local velocity field, it follows

$$u\hat{\mathbf{u}} \cdot \hat{\mathbf{n}}_{1,3} = \pm a \, ,$$

i.e.

$$\cos\theta_{1,3} = \pm\frac{a}{u} = \pm\frac{1}{M} \, .$$

For **locally subsonic flows**, $|\mathbf{u}| < a$, the eigenvalues $s_{1,3}$ are always negative and positive respectively, while the eigenvalue $s_2$ becomes equal to zero for unit vectors that are orthogonal w.r.t. the local velocity.

## Right and left eigenvectors

Right eigenvector problem reads

$$\mathbf{AR} = \mathbf{RS} \, ,$$

s.t. $\mathbf{Ar}_i = s_i\mathbf{r}_i$, with $\mathbf{r}_i$ the $i^{th}$ column of matrix $\mathbf{R}$.

Left eigenvector problem reads

$$\mathbf{LA} = \mathbf{SL} \, ,$$

s.t. $\mathbf{l}_i^T\mathbf{A} = s_i\mathbf{l}_i^T\mathbf{A}$, with $\mathbf{l}_i^T$ the $i^{th}$ row of matrix $\mathbf{L}$.

**Compatibility equations.** For **locally supersonic flows**, for $s_{1,3} = u_n \mp a$, $u_n = \pm a$ to get $s_{1,3} = 0$, and thus the left eigenvectors for those choices of unit vector become

$$\mathbf{l}_{0;1,3}^T = \frac{1}{2\rho a^2} \left[ \, a(a \pm u_n) \mid \mp a\mathbf{n}_{1,3}^T \right] =$$

$$= \frac{1}{2\rho a} \left[ \, 2a \mid \mp \mathbf{n}_{1,3}^T \right] \, .$$

The quasi linear form of the equations becomes

**todo** *Check algebra, and UNCOMMENT!*

$\ldots$

For $s_2 = u_n$, $u_n = 0$ to get $s_2 = 0$, and thus the corresponding left eigenvector reads

$$\mathbf{l}_{0;2}^T = \frac{1}{\rho a} \left[ -\mathbf{t}^T\mathbf{u} \mid \mathbf{t}^T \right] \, ,$$

with $\mathbf{t} = \begin{bmatrix} n_y \\ -n_x \end{bmatrix}$, so that $\mathbf{t}^T\mathbf{n} = 0$. The quasi linear form of the equations becomes

$$0 = \frac{1}{\rho a}\left[ -\mathbf{t}^T\mathbf{u} \mid \mathbf{t}^T \right]\left( \partial_t \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} + \begin{bmatrix} \cdot & \cdot & \cdot \\ a^2 - \frac{m_x m_x}{\rho^2} & \frac{2m_x}{\rho} & \cdot \\ -\frac{m_x m_y}{\rho^2} & \frac{m_y}{\rho} & \frac{m_x}{\rho} \end{bmatrix} \partial_x \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} + \begin{bmatrix} \cdot & \cdot & \cdot \\ -\frac{m_x m_y}{\rho^2} & \frac{m_y}{\rho} & \frac{m_x}{\rho} \\ a^2 - \frac{m_y m_y}{\rho^2} & \cdot & \frac{2m_y}{\rho} \end{bmatrix} \partial_y \begin{bmatrix} \rho \\ m_x \\ m_y \end{bmatrix} = \mathbf{0} \right) =$$

$$0 = \left[ -\mathbf{t}^T\mathbf{u} \mid \mathbf{t}^T \right]\partial_t \begin{bmatrix} \rho \\ \mathbf{m} \end{bmatrix} +$$

$$+ \left[ \, t_x a^2 - \mathbf{t}^T\mathbf{u}u_x \mid -\mathbf{t}^T\mathbf{u} + t_x u_x + \mathbf{t}^T\mathbf{u} \mid t_y u_x \, \right]\partial_x \begin{bmatrix} \rho \\ \mathbf{m} \end{bmatrix} +$$

$$+ \left[ \, t_y a^2 - \mathbf{t}^T\mathbf{u}u_y \mid t_x u_y \mid -\mathbf{t}^T\mathbf{u} + t_y u_y + \mathbf{t}^T\mathbf{u} \, \right]\partial_y \begin{bmatrix} \rho \\ \mathbf{m} \end{bmatrix} =$$

$$0 = \left[ -\mathbf{t}^T\mathbf{u} \mid \mathbf{t}^T \right]\partial_t \begin{bmatrix} \rho \\ \mathbf{m} \end{bmatrix} +$$

$$+ a^2\left( t_x\partial_x + t_y\partial_y \right)\rho - u_t\left( u_x\partial_x + u_y\partial_y \right)\rho + t_x\left( u_x\partial_x + u_y\partial_y \right)m_x + t_y\left( u_x\partial_x + u_y\partial_y \right)m_y =$$

$$=$$

or choosing $\hat{\mathbf{t}} = \hat{\mathbf{u}}$, s.t. $\mathbf{u} = u\hat{\mathbf{t}}$,

$$0 = \dots \partial_t \dots + a^2 \partial_u \rho - u^2 \partial_u \rho + u t_x \partial_u m_x + u t_y \partial_u m_y =$$
$$= \dots \partial_t \dots + a^2 \partial_u \rho - u^2 \partial_u \rho + \rho \underbrace{u t_x}_{u_x} \partial_u u_x + \rho \underbrace{u t_y}_{u_y} \partial_u u_y + \underbrace{u u_x t_x \partial_u \rho + u u_y t_y \partial_u \rho}_{= u^2 \partial_u \rho} =$$
$$= \dots \partial_t \dots + a^2 \, \partial_u \rho + \rho \partial_u \frac{|\mathbf{u}|^2}{2} \ .$$

In **steady conditions**, dividing by $\rho \neq 0$

$$0 = a^2 \frac{\partial_u \rho}{\rho} + \partial_u \frac{u^2}{2} = \partial_u \left( a^2 \ln \rho + \frac{u^2}{2} \right) \ .$$

This is the isothermal version of the Bernoulli's theorem along a streamline for compressible flows, in absence of volume force. This can be derived from momentum equation in steady conditions, after scalar multiplication by the velocity field

$$0 = \mathbf{u} \cdot [\rho \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p] =$$
$$= \mathbf{u} \cdot \left[ \rho \nabla \frac{|\mathbf{u}|^2}{2} + \omega \times \mathbf{u} + \nabla p \right] =$$
$$= \rho \mathbf{u} \cdot \left[ \nabla \frac{|\mathbf{u}|^2}{2} + a^2 \frac{\nabla \rho}{\rho} \right] =$$
$$= \rho \mathbf{u} \cdot \nabla \left( \frac{|\mathbf{u}|^2}{2} + a^2 \ln \rho \right) \ .$$

## Euler equations for inviscid compressible flows

## Shallow water

# NAVIER-CAUCHY EQUATIONS

Navier-Cauchy equations are the differential balance equation of the momentum of an elastic isotropic medium in the regime of small strain and displacement,

$$\rho_0 \partial_{tt} \vec{s} = \rho_0 \vec{g} + \nabla \cdot \boldsymbol{\sigma} .$$

Stress tensor for an isotropic medium reads

$$\boldsymbol{\sigma} = 2\mu\boldsymbol{\varepsilon} + \lambda \operatorname{tr}(\boldsymbol{\varepsilon}) \, \mathbb{I} =$$
$$= \left( 2\mu\boldsymbol{\varepsilon} - \frac{2}{3}\mu \operatorname{tr}(\boldsymbol{\varepsilon})\mathbb{I} \right) + \left( \lambda + \frac{2}{3}\mu \right) \operatorname{tr}(\boldsymbol{\varepsilon}) \, \mathbb{I} ,$$

with the small strain tensor

$$\boldsymbol{\varepsilon} = \frac{1}{2} \left( \nabla\vec{s} + \nabla^T\vec{s} \right) .$$

Essential, natural and Robin boundary conditions read

$$\vec{s} = \bar{\vec{s}} \qquad \vec{r} \in S_D \quad \text{esserntial - Dirichlet b.c.}$$
$$\hat{n} \cdot \boldsymbol{\sigma} = \bar{\vec{t}}_n \qquad \vec{r} \in S_N \quad \text{natural - Neumann b.c.}$$
$$a\vec{s} + \hat{n} \cdot \boldsymbol{\sigma} = \vec{b} \quad \vec{r} \in S_R \quad \text{Robin b.c.}$$

## 29.1 Weak formulation

For $\forall \vec{w} \in \dots$

$$0 = -\int_V \rho\vec{w} \cdot \partial_{tt}\vec{s} + \int_V \rho_0\vec{w} \cdot \vec{g} + \int_V \vec{w} \cdot \nabla \cdot \boldsymbol{\sigma} =$$
$$= -\int_V \rho\vec{w} \cdot \partial_{tt}\vec{s} + \int_V \rho_0\vec{w} \cdot \vec{g} + \int_{\partial V} \hat{n} \cdot \boldsymbol{\sigma} \cdot \vec{w} - \int_V \nabla\vec{w} : \boldsymbol{\sigma}$$

The volume integral containing the stress tensor can be written either as

$$\int_V \nabla\vec{w} : \boldsymbol{\sigma} = \int_V w_{i/j} \left[ \mu \left( s_{i/j} + s_{j/i} \right) + \lambda s_{k/k}\delta_{ij} \right] =$$
$$= \int_V \mu w_{i/j} \left( s_{i/j} + s_{j/i} \right) + \int_V \lambda w_{j/j}s_{k/k}$$

or

$$\int_V \frac{1}{2} \left( \nabla\vec{w} + \nabla^T\vec{w} \right) : \boldsymbol{\sigma} = \int_V \frac{1}{2} \left( w_{i/j} + w_{j/i} \right) \left[ \mu \left( s_{i/j} + s_{j/i} \right) + \lambda s_{k/k}\delta_{ij} \right] =$$
$$= \int_V \frac{\mu}{2} \left( w_{i/j} + w_{j/i} \right) \left( s_{i/j} + s_{j/i} \right) + \int_V \lambda w_{j/j}s_{k/k}$$

The weak formulation of the Navier-Cauchy equations reads

$$\int_V \rho_0 \vec{w} \cdot \partial_{tt} \vec{s} + \int_V 2\mu \frac{\nabla \vec{w} + \nabla^T \vec{w}}{2} : \frac{\nabla \vec{s} + \nabla^T \vec{s}}{2} + \int_V \lambda \nabla \cdot \vec{w} \nabla \cdot \vec{s} + \int_{S_R} \vec{w} \cdot a\vec{s} = \int_V \rho_0 \vec{w} \cdot \vec{g} + \int_{S_N} \vec{w} \cdot \vec{\bar{t}}_n + \int_{S_R} \vec{w} \cdot \vec{b} \,,$$

for $\forall \vec{w} \in ...$, and with $\vec{s} = \bar{\vec{s}}$ for $\vec{r} \in S_D$.

# THIRTY

# NAVIER-STOKES EQUATIONS

Incompressible Navier-Stokes equations read

$$\begin{cases} \rho \partial_t \vec{u} + \rho \left( \vec{u} \cdot \nabla \right) \vec{u} - \mu \nabla^2 \vec{u} + \nabla P = \rho \vec{g} \\ \nabla \cdot \vec{u} = 0 \, . \end{cases}$$

Mass balance equation is replaced by the incompressiblity kinematic constraint, $\nabla \cdot \vec{u} = 0$: this constraint is not dynamic, as time derivative of density does not appear in the equation. With the incompressibility constraint, mass equation tells us that material particles keep their density constant,

$$0 = \underbrace{\partial_t \rho + \vec{u} \cdot \nabla \rho}_{= \frac{D\tilde{\rho}}{Dt}} + \rho \underbrace{\nabla \cdot \vec{u}}_{=0} = \frac{D\rho}{Dt} \, ,$$

whose solution can be written using material coordinates $\vec{r}_0$ as $\rho(\vec{r}(\vec{r}_0, t), t) = \rho_0(\vec{r}_0, t)$.

## 30.1 Incompressibility constraint

Incompressibility constraint makes thermodynamic fade, while pressure field is replaced by/contains the contribution of a Lagrangian multiplier related to the incompressiblity constraint.

### 30.1.1 Wave-vector transformed space

Transforming the fields from physical space to the wave-vector space $\tilde{u}(\vec{\kappa}, t) = \mathcal{F} \{\vec{u}(\vec{r}, t)\}$, Navier-Stokes equations for incompressible fluids with uniform and constant density $\rho(\vec{r}, t) = \rho$ becomes

$$\begin{cases} \rho \partial_t \tilde{u} + \mathcal{F} \{ \left( \vec{u} \cdot \nabla \right) \vec{u} \} + \mu |\vec{\kappa}|^2 \tilde{u} + i \vec{k} \widetilde{P} = \rho \tilde{g} \\ i \vec{\kappa} \cdot \tilde{u} = 0 \, . \end{cases}$$

Taking the divergence of the momentum balance equation, i.e. taking the scalar product with $i\vec{\kappa}$ in the transformed space, and using the incompressibility constraint to set $i\vec{\kappa} \cdot \tilde{u} = 0$,

$$i\vec{\kappa} \cdot \mathcal{F} \{ \left( \vec{u} \cdot \nabla \right) \vec{u} \} - |\vec{\kappa}|^2 \widetilde{P} = i\vec{\kappa} \cdot \rho \tilde{g} \, ,$$

so that the transformed pressure field becomes

$$\widetilde{P} = \frac{i\vec{\kappa}}{|\vec{\kappa}|^2} \cdot \mathcal{F} \{ \left( \vec{u} \cdot \nabla \right) \vec{u} - \rho \vec{g} \} \, ,$$

Replacing this expression in the transformed Navier-Stokes equations, the meaning of the pressure field as a Lagrange multiplier associated with incompressibility constraint becomes clear,

$$\rho \partial_t \tilde{u} + \mu |\vec{\kappa}|^2 \tilde{u} = \left[ 1 - \frac{\vec{\kappa}\vec{\kappa}}{|\vec{\kappa}|^2} \right] \cdot \mathcal{F} \{ -(\vec{u} \cdot \nabla)\vec{u} + \rho \tilde{g} \}$$

as the orthogonal projector $[1 - \frac{\vec{\kappa}\vec{\kappa}}{|\vec{\kappa}|^2}]$ onto the space of divergence-free functions acts on the non-linear and forcing terms.

## 30.2 Weak formulation of the problem

$$0 = \int_V \vec{w} \cdot [\rho \partial_t \vec{u} + \rho(\vec{u} \cdot \nabla)\vec{u} - 2\mu \nabla \cdot \mathbb{D}(\vec{u}) + \nabla P - \rho \vec{g}] - \int_V v \nabla \cdot \vec{u} =$$

$$= \int_V \vec{w} \cdot [\rho \partial_t \vec{u} + \rho(\vec{u} \cdot \nabla)\vec{u}] + \int_V 2\mu \nabla \vec{w} : \mathbb{D} - \int_V \nabla \cdot \vec{w} P - \int_V \vec{w} \cdot \rho \vec{g} - \int_V v \nabla \cdot \vec{u} - \int_{\partial V} \hat{n} \cdot (\mathbb{S} - P\mathbb{I}) \cdot \vec{w},$$

### 30.2.1 Weak formulation and incompressibility constraint

$$\vec{r}(\vec{r}_0, t) = \vec{r}(q(t), t)$$

$$\vec{u} = \frac{D\vec{r}}{Dt} = \dot{q}\frac{\partial \vec{r}}{\partial q} + \frac{\partial \vec{r}}{\partial t}$$

In the weak formulation, using $\vec{w} = \frac{\partial \vec{r}}{\partial q} = \frac{\partial \vec{u}}{\partial \dot{q}}$

$$0 = \int_V \vec{w} \cdot \rho \frac{D\vec{u}}{Dt} + \int_V 2\mu \nabla \vec{w} : \mathbb{D} - \int_V \nabla \cdot \vec{w} P - \int_V \rho \vec{w} \cdot \vec{g} - \int_V v \nabla \cdot \vec{u} - \int_{\partial V} \vec{t}_{\hat{n}} \cdot \vec{w},$$

$$\int_V \vec{w} \cdot \rho \frac{D\vec{u}}{Dt} \, dV = \int_{V_0} \rho_0 \frac{\partial \vec{u}}{\partial \dot{q}} \cdot \frac{D\vec{u}}{Dt} =$$

$$= \int_{V_0} \rho_0 \frac{D}{Dt}\left(\frac{\partial \vec{u}}{\partial \dot{q}} \cdot \vec{u}\right) dV_0 - \int_{V_0} \rho_0 \frac{D}{Dt}\left(\frac{\partial \vec{r}}{\partial \dot{q}}\right) \cdot \vec{u} \, dV_0 =$$

$$= \int_{V_0} \rho_0 \frac{D}{Dt}\left(\frac{\partial}{\partial \dot{q}} \frac{|\vec{u}|^2}{2}\right) dV_0 - \int_{V_0} \rho_0 \frac{\partial}{\partial q} \frac{|\vec{u}|^2}{2} \, dV_0 =$$

...

## 30.3 Non-linear term

Different ways to treat the non-linear term:

- Semi-linear approximation of the non-linear term

$$(\vec{u}(\vec{r}, t^n) \cdot \nabla)\vec{u}(\vec{r}, t^n) \sim (\vec{u}^*(\vec{r}, t^n) \cdot \nabla)\vec{u}(\vec{r}, t^n),$$

with $\vec{u}^*(\vec{r}, t^n)$ an approximation of $\vec{u}(\vec{r}, t^n)$ involving values of the velocity field at previous time-steps, as an example

$$\vec{u}^*(\vec{r}, t^n) = \begin{cases} \vec{u}(\vec{r}, t^{n-1}) & 1^{st}\text{-order} \\ 2\vec{u}(\vec{r}, t^{n-1}) - \vec{u}(\vec{r}, t^{n-2}) & 2^{nd}\text{-order} \end{cases}$$

# **ARBITRARY LAGRANGIAN-EULERIAN DESCRIPTION**

Reynold's transport theorem allows for the formulation of intergal equations, and grid-based methods like FVM, on moving grids and changing domains. Rules for derivatives of composite functions provide the relations between time derivatives in a Lagrangian, Eulerian, or arbitrary description,

$$\left.\frac{\partial f}{\partial t}\right|_{\vec{r}_0} = \left.\frac{\partial f}{\partial t}\right|_{\vec{r}} + \vec{u} \cdot \nabla f$$

$$\left.\frac{\partial f}{\partial t}\right|_{\vec{r}_b} = \left.\frac{\partial f}{\partial t}\right|_{\vec{r}} + \vec{u}_b \cdot \nabla f$$

Equations governing the motion of the grid are usually required as well. E.g.:

- known and prescribed motion of the grid;

- boundary conditions only without changing grids (for small displacements)

- pseudo-elastic deformation (usually good for small strain and displacement;

- for large displacements of/or models with complex geometry, sliding and/or overlapping grids could an option for grid-based methods.

## **31.1 Integral problem**

Application of Reynolds theorem to the balance equation of the quantity $\mathbf{u}$ for a material volume $V_t$

$$\frac{d}{dt} \int_{V_t} \rho \mathbf{u} = \int_{V_t} \rho \mathbf{f} + \oint_{\partial V_t} \hat{n} \cdot \mathbf{T} \,.$$

provides the expression of the balance equation for a geometrical volume $v_t$ in arbitrary motion,

$$\frac{d}{dt} \int_{v_t} \rho \mathbf{u} + \oint_{\partial v_t} \rho \mathbf{u} \left(\vec{u} - \vec{u}_b\right) \cdot \hat{n} = \int_{v_t} \rho \mathbf{f} + \oint_{\partial v_t} \hat{n} \cdot \mathbf{T} \,.$$

Here, the integral forulation of the problem will be applied to each element of the grid in arbitrary motion, for domains with variable geometry.

## 31.2 Differential problem

Rules for derivatives of composite functions allows to write the differential w.r.t. the variables associated with the points of a moving grid. A balance equation in convective form can be written as

$$\rho \frac{D\mathbf{u}}{Dt} = \rho \mathbf{f} + \nabla \cdot \mathbf{T}$$

$$\rho \left[ \frac{\partial \mathbf{u}}{\partial t} + \vec{u} \cdot \nabla \mathbf{u} \right] =$$

$$\rho \left[ \left. \frac{\partial \mathbf{u}}{\partial t} \right|_{\vec{r}_b} + (\vec{u} - \vec{u}_b) \cdot \nabla \mathbf{u} \right] =$$

**Part X**

# Numerical Methods for PDEs

# INTRODUCTION TO NUMERICAL METHODS FOR PDES

Different numerical methods for PDEs rely on the discretization of different formulations of the continunuos problems. As an example,

- **FDM, Finite difference methods** rely on the approximation of derivatives of the **strong formulation** of the problem

- *FEM, Finite element methods* rely on a finite dimensional approximation of the **weak formulation** of the problem; usually the finite dimensional approximation can be interpreted as a projection of a infinite dimensional continuous problem onto a finite dimensional space, the space of the choosen finite elements

- *FVM, Finite volume methods* rely on an approximation of the **integral formulation** of the problem

- **BEM, Boundary elemement methods** rely on an approximation of a **boundary integral formulation** of the problem, when it's feasible and convenient

- …*spectral methods*, *spectral element methods*,…

**Characteristics.**

- grid: domain-grid-based, buondary-grid-based, grid-free methods

- range of interaction: short in physical space for FDV, FEM, FVM; long-range for boundary element methods, even though clustering techniques are available, like FMM; (usually) over the whole domain in space, short-range interaction in wave-number space for spectral methods;

**Pros and cons.** *More suited methods for each problems…; domain, order,…*

Let's take Poisson equations for a scalar function $u(\vec{r})$ to show all the possible approaches above. Here we start from the most general version of the equation, namely the integral form. As it's shown in Continuum Mechanics:Governing Equations, the most general form of balance equations is the integral form, while the differential form can be seamlessly derived only when the quantities involved are regular enough, to apply Stokes' theorem and for the derivatives appearing to exist.

The problem of interest can be interpreted as the problem of finding the temperature field $u(\vec{r})$ during steady heat conduction in the domain $\Omega$, with distributed volume heat source $f(\vec{r})$. Fourier's law assume that condcution heat flux is proportional to the gradient of the temperature, $\vec{q} = -k\nabla u$. Temperature is prescribed, $u = g$ on the region of the boundary $S_D$, while the heat flux is prescribed $\hat{n} \cdot \vec{q} = h$ on the region of the boundary $S_N$.

The **integral formulation** of the problem for all the boundary reads

$$0 = -\oint_{\partial\Omega} \hat{n} \cdot \vec{q} + \int_{\Omega} f \ . \tag{32.1}$$

**Example 32.1 (Finite volume methods)**

Finite volume methods rely on a tassellation of the domain $\Omega = \cup_k \Omega_k$, to write and solve the integral balance equation (32.1) for all the elemantary domain $\Omega_k$.

$$0 = -\oint_{\partial\Omega_k} \hat{n} \cdot \vec{q} + \int_{\Omega_k} f \; . \tag{32.2}$$

evaluating volume integrals with the internal variables of the cell $\Omega_k$, and boundary flux with the variables of the cell $\Omega_k$ and the neighbouring cells $\Omega_i \in B_k$. Introducing the definition of numerical flux $F_{ik}(u_i, u_k)$ at the interface between a generic discrete version of the balance equation (32.2) for the $k^{th}$ element becomes

$$0 = -\sum_{\Omega_i \in B_k} F_{ik}(u_i, u_k) + S_k(u_k) \; .$$

Different numerical methods differs in the evaluation of fluxes and sources. FVM is **conservative** as it evaluate flux at interfaces and then distribute it to neighoring cells, see *Property 34.1*. A rough elementwise-uniform variable with jump at interfaces and diffusive flux,

$$F_{ik} = -A_{ik}\, k\, \frac{u_i - u_k}{d_{ik}} \; ,$$

with $d_{ik}$ equal to the distance of the centers of elements $i$ and $k$, leads to the balance equation for the internal volume $\Omega_k$,

$$0 = \sum_i A_{ik}\, k\, \frac{u_i - u_k}{d_{ik}} + V_k f_k \; .$$

Volumes at the boundary are influenced by fluxes at the boundaries and boundary conditions in general.

For regular cubic mesh, $A_{ik} = \Delta x$, $d_{ik} = \Delta x$, $V_k = \Delta x$, and thus the equations for internal elements become

$$0 = k \sum_{i \in \Omega_k} \frac{u_i - u_k}{\Delta x^2} + f_k \; ,$$

where the summation is exactly equal to the center-difference stencil for the Laplacian, used in finite difference method *Example 32.2*.

**If fields are regular enough** to apply Stokes' theorem, it's possible to derive differential problem from the integral formulation,

$$0 = -\oint_{\partial V} \hat{n} \cdot \vec{q} + \int_V f = -\int_V \nabla \cdot \vec{q} + \int_V f = \int_V (-\nabla \cdot \vec{q} + f)$$

Exploiting the arbitrariety of the volume $V$, the **strong form** of the differential problem is the Poisson equation with suitable boundary conditions,

$$\begin{cases} -\nabla \cdot (k\nabla u) = f(\vec{r}) \; , & \vec{r} \in V \\ u = g(\vec{r}) \; , & \vec{r} \in S_D \\ \hat{n} \cdot \nabla u = h(\vec{r}) \; , & \vec{r} \in S_N \end{cases}$$

only holds in regions of the domain where the functions involved are continuous and differentiable, i.e. everything written in the problem is not meaningful, i.e. at least it exists.

**Example 32.2 (Finite difference methods)**

Finite different method approximates the strong form of the problem, building a stencil to evaluate an approximation of

the Laplacian. On a cubic regular grid ($\Delta x = \Delta y = \Delta z$), the Laplacian van be evaluated with a 7-point stencil

$$\nabla^2 u = \partial_{xx} u + \partial_{yy} u + \partial_{zz} u =$$
$$= \frac{u_{i+1,j,k} - 2u_{i,j,k} + u_{i-1,j,k}}{\Delta x^2} + \frac{u_{i,j+1,k} - 2u_{i,j,k} + u_{i,j-1,k}}{\Delta y^2} + \frac{u_{i,j,k+1} - 2u_{i,j,k} + u_{i,j,k-1}}{\Delta z^2} =$$
$$= \frac{1}{\Delta x^2} \left[ -6u_{i,j,k} + u_{i+1,j,k} + u_{i-1,j,k} + u_{i,j+1,k} + u_{i,j-1,k} + u_{i,j,k+1} + u_{i,j,k-1} \right] =$$
$$= \frac{1}{\Delta x^2} \sum_{i \in B_k} (u_i - u_k)$$

and it's equal to the same expression already provided for the rough finite volume method with regular grid, in *Example 32.1*

Starting from strong formulation of the differential problem, the **weak formulation** of the problem is derived:

- multiplying the strong problem by an arbitrary test function $w(\vec{r})$ compatible with the essential constraints

- integrating over the whole domain

$$0 = \int_V w(\vec{r}) \cdot [-\nabla \cdot (k\nabla u) - f(\vec{r})] =$$
$$= \int_V \{ k\nabla w(\vec{r}) \cdot \nabla u(\vec{r}) - w(\vec{r})f(\vec{r}) \} - \oint_{\partial V} w(\vec{r})\,\hat{n} \cdot \nabla u \qquad \forall w(\vec{r}) .$$

Using a test function $w(\vec{r})$ equal to zero on the boundary where essential boundary conditions are prescribed $w|_{S_D} = 0$, and introducing the natural boundary conditions of the Neumann boundary $S_N$, $\hat{n} \cdot \nabla u|_{S_N} = h$,

$$\int_V k\nabla w \cdot \nabla u = \int_V wf + \int_{S_N} wh , \quad \forall w$$

compatible with essential boundary conditions.

---

**Example 32.3 (Finite element methods)**

Finite element methods build an $N$-dimensional systems:

- approximating the solution as a linear combination of $N$ base functions, $u(\vec{r}) = \sum_j \phi_j(\vec{r})u_j$

- testing the equation over $N$ independent test functions $\psi_i(\vec{r})$

i.e. the $i^{th}$ equation becomes

$$\int_V k\nabla\psi_i(\vec{r}) \cdot \nabla\phi_j(\vec{r})\,u_j = \int_V \psi_i(\vec{r})f + \int_{S_N} \psi_i(\vec{r})h ,$$

or briefly

$$K_{ij}u_j = f_i \qquad , \qquad \mathbf{Ku} = \mathbf{f} .$$

A common choice uses the same functions both as test and base functions, $\psi_i(\vec{r}) = \phi_i(\vec{r})$.

---

The differential problem in strong form can be recast as a boundary element problem exploiting the properties of Green's functions $G(\vec{r}; \vec{r}_0)$. The expression of *Green's function $G(\vec{r}, \vec{r}_0)$* is known for some problems of interest like Poisson

equation, Helmholtz equation or wave equation. Exploiting the properties of Green's function and integration by parts, the integral buondary problem is derived as follow

$$
\begin{aligned}
E(\vec{r}_0)u(\vec{r}_0) = \int_{\vec{r} \in V} u(\vec{r})\delta(\vec{r} - \vec{r}_0)\, dV = \\
= -\int_{\vec{r} \in V} u(\vec{r})\nabla^2 G(\vec{r}; \vec{r}_0)\, dV = \\
= -\oint_{\vec{r} \in \partial V} u(\vec{r})\hat{n}(\vec{r}) \cdot \nabla G(\vec{r}; \vec{r}_0)\, dS + \oint_{\vec{r} \in \partial V} G(\vec{r}; \vec{r}_0)\hat{n}(\vec{r}) \cdot \nabla u(\vec{r})\, dS - \int_{\vec{r} \in V} G(\vec{r}; \vec{r}_0)\nabla^2 u(\vec{r})\, dV = \\
= -\oint_{\vec{r} \in \partial V} u(\vec{r})\hat{n}(\vec{r}) \cdot \nabla G(\vec{r}; \vec{r}_0)\, dS + \oint_{\vec{r} \in \partial V} G(\vec{r}; \vec{r}_0)\hat{n}(\vec{r}) \cdot \nabla u(\vec{r})\, dS + \int_{\vec{r} \in V} G(\vec{r}; \vec{r}_0)f(\vec{r})\, dV\;.
\end{aligned}
$$

It must be paid attention that the integrals may be singular and may produce discontintuities in the fields across the surface $\partial V$ (*so what's the value of a field in presence of a discontinuity?...virtual singularities, regularization close to the singularities...one way to correctly interpret them is via Cauchy principal value... evaluating some integrals in a check-point just inside or outside the domain may make the value of the function $E(\vec{r}_0)$ change, but the integral involving $\hat{n} \cdot \nabla G$ - related to the solid angle seen by the check-point of the surface - changes accordingly and keep the sum of these two terms constant=*.

The value of the unknown function or the flux $\hat{n} \cdot \nabla u$ are known on the Dirichlet and Neumann regions of the boundary respectively, and so the integro-differential problem becomes

$$
E(\vec{r}_0)u(\vec{r}_0) + \int_{S_N} u\hat{n} \cdot \nabla G\, dS - \int_{S_D} G\,\hat{n} \cdot \nabla u\, dS = \int_{S_D} g\hat{n} \cdot \nabla G\, dS - \int_{S_N} Gh\, dS + \int_{\vec{r} \in V} Gf\, dV
$$

The unknown function is approximated on the boundary of the domain $\partial V$ as an $N$-dimensional approximation, as an example

$$
u(\vec{r}) = \sum_j \phi_j(\vec{r})u_j \qquad , \qquad \vec{r} \in \partial V\;,
$$

and the integro-differential equation is evaluated in $N$ different points $\vec{r}_{0,i}$ in order to get a $N, N$ linear equation in the amplitudes $u_j$ of the base functions,

$$
[\mathbf{E} + \mathbf{D}_N + \mathbf{S}_D]\,\mathbf{u} = \mathbf{D}_D\mathbf{g} + \mathbf{S}_N\mathbf{h} + \mathbf{S}_V\mathbf{f}\;.
$$

# FINITE ELEMENT METHOD

## Examples

- Structural mechanics of linear beam structures

- *Poisson equation, 1D*

## 33.1 1-dimensional Poisson equation

### 33.1.1 Strong form of the problem

$$- \left( \nu u'(x) \right)' = f(x) \quad , \qquad x \in [a, b]$$

Supplied with proper boundary conditions. As an example,

- Dirichlet boundary conditions: $u(\overline{x}) = \overline{u}$

- Neumann boundary conditions: $\nu u'(\overline{x}) = \overline{q}$

- Robin boundary conditions: $au(\overline{x}) + bu'(\overline{x}) = c$

### 33.1.2 Weak form of the problem

For $\forall w(x) \in \mathcal{H}^1([a, b])$[1] ,

$$0 = \int_{x=a}^{b} w(x) \left\{ - \left( \nu u'(x) \right)' - f(x) \right\} =$$

$$= \int_{x=a}^{b} w'(x) \nu u'(x) \, dx - \int_{x=a}^{b} w(x) f(x) \, dx - \left[ w(x) \nu u(x) \right]\big|_{x=a}^{b}$$

Boundary conditions…

---

[1] Which functional space? It would be nice to be as more precise as possible, without introducing too many "unnecessary" complications. Without going into details, 1) everything appears in the weak form should exists for the functions of that space: here, the functions must have $L^2$-integrable first-order derivative; 2) some results exist when test and unkwown functions belong to the same space: so, some "special" treatment of essential boundary conditions (on Dirichlet boundaries) is required.

## Poisson equation with Dirichlet boundary conditions

By definition, test functions are identically zero on Dirichlet boundary conditions $w|_{S_D} = 0$.

Thus the weak form becomes

$$\int_{x=a}^{b} w'(x)\nu u'(x)\,dx = \int_{x=a}^{b} w(x)f(x)\,dx \,,$$

along with the essential boundary conditions $u(a) = u_a$, $u(b) = u_b$.

## Finite element method

**Discretization of the domain.** The domain $[a, b] =: \Omega$, is divided into elements $\Omega_i$ s.t.

$$\cup_i \Omega_i = \Omega$$
$$\Omega_i \cap \Omega_j = \emptyset \quad i \neq j$$

Exploiting additive property of integrals on the union of domains, integrals of the weak form become summation of integrals on the individual elements

$$\int_\Omega \cdots = \sum_i \int_{\Omega_i} \cdots$$

**Choice of the base functions.** *Lagrangian base functions* (1) $\phi_i(x_j) = \delta_{ij}$, 2) $\phi_i(x) = 0$, if $x \notin B(x_i)$; 3) $\sum_i \phi_i(x) = 1, \forall x \in \Omega$)

$$u(x) = \sum_i \phi_i(x)u_i \,.$$

Definition on a reference element $\xi \in [-1, 1]$ and then change of coordinates between "physical" and reference space to evaluate integrals.

Here, as the maximum order of the derivatives involved in the problem is 1, test functions must be piecewise-*linear* (or better, first degree polynomial) at least, in order to avoid numerically setting to zero some terms in the equation only due to poor approximation (here the term $\int_{x=a}^{b} w'(x)\nu u'(x)dx$ would be zero with piecewise-constant test functions).

First degree-polynomial Lagrangian base functions on the reference elements are

$$\widetilde{\varphi}_1(\xi) = \frac{1}{2}(1 - \xi)$$
$$\widetilde{\varphi}_2(\xi) = \frac{1}{2}(1 + \xi)$$

**Evaluation of integrals.** Polynomials functions are exactly evaluated using Gaussian integration, i.e. with a linear combination of the values of the function on a set of Gaussian nodes $x\_g$

…

**Analytical exact integration for first degree-polynomial.** The coordinate transformation[2]

$$x = a_i \frac{1 - \xi}{2} + b_i \frac{1 + \xi}{2} = \frac{a_i + b_i}{2} + \xi \frac{\ell_i}{2} \,,$$

transforms the reference element into physical elements: e.g. $x(\xi = -1) = a_i$, $x(\xi = 1) = b_i$. The Jacobian of this transformation is constant over the element,

$$\frac{dx}{d\xi} = \frac{\ell_i}{2} \,.$$

---

[2] Here, the transformation between physical and reference space is iso-parametric, as it uses the same base functions as those used for function approximations, $x(\xi) = a_i \widetilde{\phi}_1(\xi) + b_i \widetilde{\phi}_2(\xi)$.

---

On the element $i$,

$$K_{i,11} = \int_{x=a_i}^{b_i} \frac{d\phi_1}{dx}(x)\frac{d\phi_1}{dx}(x)\, dx =$$

$$= \int_{\xi=-1}^{\xi=1} \frac{d\xi}{dx}\frac{d\tilde{\phi}_1}{d\xi}\frac{d\xi}{dx}\frac{d\tilde{\phi}_1}{d\xi}\frac{dx}{d\xi}d\xi =$$

$$= \int_{\xi=-1}^{\xi=1} \left(-\frac{1}{2}\right)\left(-\frac{1}{2}\right)d\xi\frac{2}{\ell_i} = \frac{1}{\ell_i} \;.$$

Analogously,

$$K_{i,11} = K_{i,22} = \frac{1}{\ell_i}$$

$$K_{i,12} = K_{i,21} = -\frac{1}{\ell_i}$$

$$M_{i,11} = \int_{x=a_i}^{b_i} \phi_1(x)\phi_1(x)dx =$$

$$= \int_{\xi=-1}^{1} \varphi_1(\xi)\varphi_1(\xi)\frac{dx}{d\xi}d\xi =$$

$$= \int_{\xi=-1}^{1} \left[\frac{1}{2}(1-\xi)\right]^2 d\xi\frac{\ell_i}{2} =$$

$$= \frac{\ell_i}{8}\left(2+\frac{2}{3}\right) = \frac{1}{3}\ell_i \;.$$

and analogously

$$M_{i,11} = M_{i,22} = \frac{\ell_i}{3}$$

$$M_{i,12} = M_{i,21} = \frac{\ell_i}{6}$$

```python
import numpy as np
import scipy as sp

import matplotlib.pyplot as plt

#> Domain
a, b = 0, 1
l = b - a

#> Physical properties
nu = .01                        # Diffusion coefficient
f = 1.

#> Discretization
nel = 20            # n.of elements
nnodes = nel + 1
#> Dirichlet boundaries
i_dir = np.array([ 0, nel ])    # Global idx of nodes on Dirichlet boundaries
u_dir = np.array([ .0, .0 ])    # Value of the solution at nodes on Dir. boundaries

#> Volume forcing contribution
```

(continues on next page)

```
f_nodal = np.ones(nel+1) * f    # Nodal values f(x_i) of volume forcing, here uniform␣
 ↪f(x) = f

#> Node coordinates array, rr, and element-node connectivity array, ee
rr = np.linspace(a,b, nel+1)
ee = np.array([[i, i+1] for i in np.arange(nel)])
#> Size of the elements
e_vol = np.array([ rr[ee[i,1]] - rr[ee[i,0]] for i in np.arange(nel) ])
```

```
#> Assembling stiffenss matrix, mass matrix, and forcing term
#> Stiffness matrix
ki = np.concatenate([[ee[i,0], ee[i,0], ee[i,1], ee[i,1]] for i in np.arange(nel)])
kj = np.concatenate([[ee[i,0], ee[i,1], ee[i,0], ee[i,1]] for i in np.arange(nel)])
ke = np.concatenate([nu * np.array([1., -1., -1., 1.,])/e_vol[i] for i in np.
 ↪arange(nel)])
K = sp.sparse.coo_array((ke, (ki, kj)),)

#> Mass matrix, here used for integration of the volume force
mi = np.concatenate([[ee[i,0], ee[i,0], ee[i,1], ee[i,1]] for i in np.arange(nel)])
mj = np.concatenate([[ee[i,0], ee[i,1], ee[i,0], ee[i,1]] for i in np.arange(nel)])
me = np.concatenate([np.array([2., 1., 1., 2.,])*e_vol[i]/6. for i in np.arange(nel)])
M = sp.sparse.coo_array((me, (mi, mj)),)

F = M @ f_nodal    # Volume force
```

### Slicing the linear system

```
#> Apply essential boundary conditions, solve the linear system and retrieve the␣
 ↪solution
#> Method 1. Slicing matrices
# K u = f
# [ Kuu  Kud ] [ uu ] = [ fu ]
# [ Kdu  Kdd ] [ ud ]   [ fd ]
# with known: ud, fu; unknown: uu, fd found solving the problem
# Kuu * uu = fuu - Kud * ud  -> uu = ...
# fd = Kdu * uu + Kdd * ud
iD = i_dir.copy()
iU = np.array(list(set(np.arange(nnodes)) - set(iD)))
ud = u_dir.copy()

#> Slice stiffness matrix
Kuu = (K.tocsc()[:,iU]).tocsr()[iU,:]
Kud = (K.tocsc()[:,iD]).tocsr()[iU,:]
Kdu = (K.tocsc()[:,iU]).tocsr()[iD,:]
Kdd = (K.tocsc()[:,iD]).tocsr()[iD,:]

#> Slice forcing
fu = F[iU]

#> Solve the linear system
uu = sp.sparse.linalg.spsolve(Kuu, fu - Kud @ ud )
fd = Kdu @ uu + Kdd @ ud
```
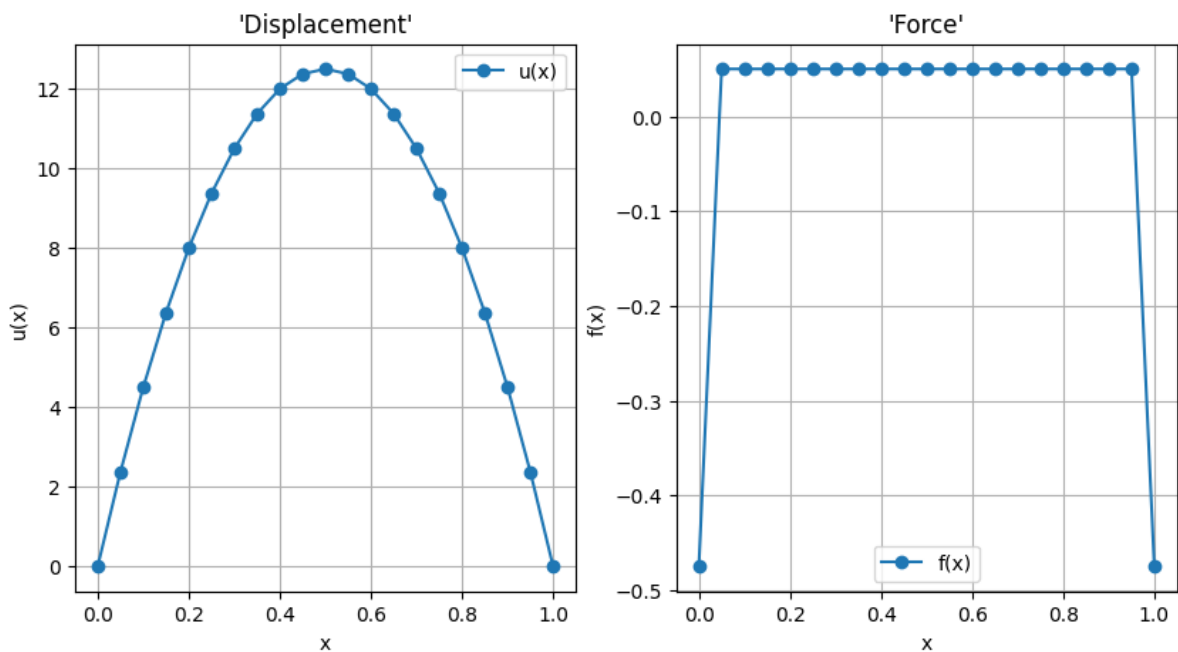
```python
#> Re-assemble the solution
u = np.zeros(nnodes);  u[iU] = uu;   u[iD] = ud
f = np.zeros(nnodes);  f[iU] = fu;   f[iD] = fd
# print("uu: \n", uu)
# print("Fd: \n", fd)

#> Plot some results
fig, ax = plt.subplots(1,2, figsize=(10, 5))
ax[0].plot(rr, u, '-o', label='u(x)')
ax[0].legend()
ax[0].grid()
ax[0].set_title('\'Displacement\'')
ax[0].set_xlabel('x')
ax[0].set_ylabel('u(x)')
ax[1].plot(rr, f, '-o', label='f(x)')
ax[1].legend()
ax[1].grid()
ax[1].set_title('\'Force\'')
ax[1].set_xlabel('x')
ax[1].set_ylabel('f(x)')
```

```
Text(0, 0.5, 'f(x)')
```



```python
#> Method 2. Augmented system
# ...
```

## Augmenting the linear system

See discussion at the end of the *script implementing a basic finite volume method solution of 1-dimensional Poisson equation*.

**Some comments and todos.**

- **Convergence analysis.** It's possible to evaluate both convergence with the dimension of the elements and/or the degree of the polynomial base functions. The *exact solution* can be easily computed analytically via direct integration for simple distribution of load $f(x)$, and uniform $\nu$

- **Result discussion.** "Force" contribution on Dirichlet boundary conditions involves both the node contribution to distributed load and the "constraint reaction". As an example, with a domain of size $\ell$ and uniform force distribution $f$, with $u(0) = u(\ell) = 0$, by *symmetry* and *global equilibrium* the force at boundaries are $-\frac{1}{2}q\ell = -\frac{1}{2}$. With the modelling choices done, the contribution of the uniform distributed load on an element on any of its nodes is $\frac{1}{2}q\ell_i = \frac{1}{2}q\frac{\ell}{n_{el}}$. Here with $n_{el} = 20$, this elementary contribution is $\frac{1}{40} = .025$: thus, the nodal contribution of internal nodes is twice (contribution from two neighboring elements) this value .05, and the force on each Dirichlet boundary is $-\frac{1}{2} + \frac{1}{40} = -0.475$.

- Try to set different values of Dirichlet boundary conditions, u_dir

# FINITE VOLUME METHOD

---

**Property 34.1**

Evaluate flux on interfaces between cells, and distribute between neighboring cells.

---

## 34.1 1-dimensional Poisson equation

### 34.1.1 Integral form of the problem

The integral form of a Poisson problem reads, **for every** $V \in \Omega$,

$$- \oint_{\partial V} \hat{\mathbf{n}} \cdot (\nu \nabla u) = \int_V f \, .$$

supplied with proper boundary conditions. As an example,

- Dirichlet boundary conditions: $u(\overline{x}) = \overline{u}$
- Neumann boundary conditions: $\nu \hat{\mathbf{n}} \cdot \nabla u(\overline{x}) = \overline{q}$
- Robin boundary conditions: $au(\overline{x}) + bu'(\overline{x}) = c$

*The multi-dimensional set may be useful also for 1-dimensional domains, to avoid confusions with signs*

### 34.1.2 Domain, volume forcing and boundary conditions

```python
import numpy as np
import scipy as sp

import matplotlib.pyplot as plt

#> Domain
a, b = 0, 1
l = b - a

#> Physical properties
nu = .01                        # Diffusion coefficient
f = 1.

#> Discretization
```

```python
nel = 10              # n.of elements
nnodes = nel + 1      # n.of nodes
nvol = nnodes         # n.finite volumes (node-centered FVM)

#> Dirichlet boundaries
i_dir = np.array([ 0, nvol-1])   # Global idx of nodes on Dirichlet boundaries
u_dir = np.array([ .0, .0 ])   # Value of the solution at nodes on Dir. boundaries

#> Volume forcing contribution
f_nodal = np.ones(nel+1) * f   # Nodal values f(x_i) of volume forcing, here uniform␣
↪f(x) = f

#> Node centered FVM
rr_n = np.linspace(a,b, nel+1)       # node coordinates (center of the node-centered␣
↪FVs)
rr_i = .5 * ( rr_n[:-1] + rr_n[1:] ) # coords of the inner cell interface
rr_b = np.array([ 2*rr_n[0]-rr_n[1], 2*rr_n[-1]-rr_n[-2] ])   # coords of the␣
↪boundaries

vols = np.array(
    [ np.abs(rr_i[0]-rr_n[0]) ] + \
      list( np.abs(rr_i[1:] - rr_i[:-1]) ) + \
    [ np.abs(rr_i[-1]-rr_n[-1]) ]
)

ni = len(rr_i)
nb = len(rr_b)

#> Interface-volume connectivity (ni, 2) (normal dir: first element -> second element)
iee = np.array([ [i, i+1] for i in np.arange(ni) ])
#> Boundary-volume connectivity (nb) (convention, normal dir: elem -> outside)
bee = np.array([0, nvol-1])
```

### 34.1.3 Discrete matheamtical problem

```python
#> Evaluating boundary contributions
# Here assmbling a stiffness matrix. In many FVM applications, there's explicit time
# dependence and explicit time integration schemes are used, so there's no need to
# assemble any matrix while the algorithm just rely on accumulation of flux␣
↪contributions
# -> Add a link to a script implementing time-dependent "solver" for FVM, e.g.␣
↪hyperbolic problems

i_flux, j_flux, e_flux = [], [], []

#> Loop over internal interface
# accumulating and keep changing the dimension! May be very inefficient and slow
for i in np.arange(ni):
    flux_u = nu / np.abs( rr_n[ iee[i,1] ] - rr_n[ iee[i,0] ] )
    i_flux += [ iee[i,0], iee[i,0], iee[i,1], iee[i,1] ]
    j_flux += [ iee[i,0], iee[i,1], iee[i,0], iee[i,1] ]
    e_flux += [ -flux_u, flux_u, flux_u, -flux_u ]

e_flux = -np.array(e_flux)
```

```
K = sp.sparse.coo_array((e_flux, (i_flux, j_flux)),)

#> Volume forcing
F = vols * f

#> Loop over boundary surfaces
# ...
# - Neumann: add a forcing term for finite volumes with sides on Neumann boundaries
# - Dirichlet: prescribe node value (e.g. 1) matrix slicing, 2) augmented system,...)
# - Robin: add a forcing term depending on the unknown value of the function ->␣
 ↪modify K matrix

# ...
```

### 34.1.4 Applying essential boundary conditions

The discrete counterpart of the continuous problem is a linear system that can be formally written as $\mathbf{Au} = \mathbf{f}$. So far, the essential boundary conditions on Dirichlet boundary $S_D$ have not been prescribed yet. Essential boundary conditions can be prescribed in (at least) two ways: 1) slicing the linear system; 2) augmenting the linear system.

**Method 1. Slicing the linear system.** The set of nodes $(\cdot)_d$ lying on Dirichlet boundaries can be distinguished from all the remaining nodes $(\cdot)_u$, and the linear system sliced accordingly

$$\begin{bmatrix} \mathbf{A}_{uu} & \mathbf{A}_{ud} \\ \mathbf{A}_{du} & \mathbf{A}_{dd} \end{bmatrix} \begin{bmatrix} \mathbf{u}_u \\ \mathbf{u}_d \end{bmatrix} = \begin{bmatrix} \mathbf{f}_u \\ \mathbf{f}_d \end{bmatrix} \ .$$

On $u$-nodes, forcing $\mathbf{f}_u$ is known, while the value of the function $u(\mathbf{x})$ is unknown, and thus $\mathbf{u}_n$ is unknown; on the other hand, on $d$-nodes the value of the function $u(\mathbf{x})$ is known, and so $\mathbf{u}_d = \mathbf{u}_D$ is known, while the forcing $\mathbf{f}_d$ is not known (as it contains both a volume contribution and a "boundary contribution" required to prescribe essential boundary conditions, as it will be more clear later, hopefully).

After slicing, the problem is solved first finding $\mathbf{u}_u$ and then retrieving the value of $\mathbf{f}_d$,

$$\begin{aligned} \mathbf{A}_{uu}\mathbf{u} = \mathbf{f}_u - \mathbf{A}_{ud}\mathbf{u}_D & \qquad \rightarrow \qquad \mathbf{u}_u = \mathbf{A}_{uu}^{-1}\left(\mathbf{f}_u - \mathbf{A}_{ud}\mathbf{u}_D\right) \\ \mathbf{f}_d = \mathbf{A}_{du}\mathbf{u}_u + \mathbf{A}_{dd}\mathbf{u}_D & \qquad \rightarrow \qquad \mathbf{f}_d = ... \end{aligned}$$

**Method 2. Augmenting the linear system.** The linear system is augmented, a) explicityl adding equations $\mathbf{u}_d = \mathbf{u}_D$ for the essential boudary conditions, and splitting the forcing contribution on Dirichlet nodes as the sum of a known volume contribution and an unknown boundary contribution,

$$\mathbf{f} = \begin{bmatrix} \mathbf{f}_u^{vol} \\ \mathbf{f}_d^{vol} + \mathbf{f}_d^D \end{bmatrix} = \begin{bmatrix} \mathbf{f}_u^{vol} \\ \mathbf{f}_d^{vol} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix} \mathbf{f}_d^D \ .$$

Moving all the unkowns on the LHS, the augmented linear system reads

$$\begin{bmatrix} \mathbf{A}_{uu} & \mathbf{A}_{ud} & \mathbf{0} \\ \mathbf{A}_{du} & \mathbf{A}_{dd} & -\mathbf{I} \\ \mathbf{0} & -\mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_u \\ \mathbf{u}_d \\ \mathbf{f}_d^D \end{bmatrix} = \begin{bmatrix} \mathbf{f}_u^{vol} \\ \mathbf{f}_d^{vol} \\ -\mathbf{u}_D \end{bmatrix} \ .$$

The (arbitrary) choice of the sign in the last block of equations (the ones prescribing the value of the unknown function on Dirichlet boundary) is usually chosen to preserve the symmetry of the original problem, if any.

**todo** *Discuss the reason why the system becomes determinate, after the application of the essential b.c.s*

### Slicing the linear system

Slice the linear system to prescribe the essential boundary conditions, solve the linear system, retrieve and plot the results.

```python
#> Apply essential boundary conditions, solve the linear system and retrieve the
 ↪solution
#> Method 1. Slicing matrices
# K u = f
# [ Kuu  Kud ] [ uu ] = [ fu ]
# [ Kdu  Kdd ] [ ud ]   [ fd ]
# with known: ud, fu; unknown: uu, fd found solving the problem
# Kuu * uu = fuu - Kud * ud  -> uu = ...
# fd = Kdu * uu + Kdd * ud
iD = i_dir.copy()
iU = np.array(list(set(np.arange(nnodes)) - set(iD)))
ud = u_dir.copy()

print(iD)
print(iU)

#> Slice stiffness matrix
Kuu = (K.tocsc()[:,iU]).tocsr()[iU,:]
Kud = (K.tocsc()[:,iD]).tocsr()[iU,:]
Kdu = (K.tocsc()[:,iU]).tocsr()[iD,:]
Kdd = (K.tocsc()[:,iD]).tocsr()[iD,:]

#> Slice forcing
fu = F[iU]

#> Solve the linear system
uu = sp.sparse.linalg.spsolve(Kuu, fu - Kud @ ud )
fd = Kdu @ uu + Kdd @ ud

#> Re-assemble the solution
u = np.zeros(nnodes);  u[iU] = uu;  u[iD] = ud
f = np.zeros(nnodes);  f[iU] = fu;  f[iD] = fd
# print("uu: \n", uu)
# print("Fd: \n", fd)

#> Plot some results
fig, ax = plt.subplots(1,2, figsize=(10, 5))
ax[0].plot(rr_n, u, '-o', label='u(x)')
ax[0].legend()
ax[0].grid()
ax[0].set_title('\'Displacement\'')
ax[0].set_xlabel('x')
ax[0].set_ylabel('u(x)')
ax[1].plot(rr_n, f, '-o', label='f(x)')
ax[1].legend()
ax[1].grid()
ax[1].set_title('\'Force\'')
ax[1].set_xlabel('x')
ax[1].set_ylabel('f(x)')
```

```
[ 0 10]
[1 2 3 4 5 6 7 8 9]
```

```
Text(0, 0.5, 'f(x)')
```



### Augmenting the linear system

Augmenting the linear system to prescribe the essential boundary conditions, solve the linear system, retrieve and plot the results.

**todo**

```
#> Method 2. Augmented system
# ...
```

**Some comments and todos.**

- **Convergence analysis.** It's possible to evaluate both convergence with the dimension of the elements and/or the degree of the polynomial base functions. The *exact solution* can be easily computed analytically via direct integration for simple distribution of load $f(x)$, and uniform $\nu$

- **Result discussion.** "Force" contribution on Dirichlet boundary conditions involves both the node contribution to distributed load and the "constraint reaction". As an example, with a domain of size $\ell$ and uniform force distribution $f$, with $u(0) = u(\ell) = 0$, by *symmetry* and *global equilibrium* the force at boundaries are $-\frac{1}{2}q\ell = -\frac{1}{2}$. With the modelling choices done, the contribution of the uniform distributed load on an element on any of its nodes is $\frac{1}{2}q\ell_i = \frac{1}{2}q\frac{\ell}{n_{el}}$. Here with $n_{el} = 20$, this elementary contribution is $\frac{1}{40} = .025$: thus, the nodal contribution of internal nodes is twice (contribution from two neighboring elements) this value .05, and the force on each Dirichlet boundary is $-\frac{1}{2} + \frac{1}{40} = -0.475$.

- Try to set different values of Dirichlet boundary conditions, u_dir

**34.1. 1-dimensional Poisson equation**                                                                 **181**

# BOUNDARY ELEMENT METHOD

Chapter 35.  Boundary Element Method

# Part XI

# Boundary Methods for PDEs

# THIRTYSIX

# GREEN'S FUNCTION METHOD

## 36.1 Poisson equation

General Poisson's problem

$$\begin{cases} -\nabla^2 \mathbf{u}(\mathbf{r}, t) = \mathbf{f}(\mathbf{r}, t) \\ + \text{ b.c.} \end{cases}$$

with common boundary conditions

$$\begin{cases} \mathbf{u} = \mathbf{g} & \text{on } S_D \\ \hat{\mathbf{n}} \cdot \nabla \mathbf{u} = \mathbf{h} & \text{on } S_N \end{cases}$$

over Dirichlet and Neumann regions of the boundary.

Poisson's problem for Green's function, in infinite domain

$$-\nabla_{\mathbf{r}}^2 G(\mathbf{r}; \mathbf{r}_0) = \delta(\mathbf{r} - \mathbf{r}_0)$$

Green's function method

$$E(\mathbf{r}_0, t) u_i(\mathbf{r}_0, t) = \int_{\mathbf{r} \in \Omega} u_i(\mathbf{r}, t) \delta(\mathbf{r} - \mathbf{r}_0) =$$

$$= -\int_{\mathbf{r} \in \Omega} u_i(\mathbf{r}, t) \nabla_{\mathbf{r}}^2 G(\mathbf{r} - \mathbf{r}_0) =$$

$$= -\int_{\mathbf{r} \in \Omega} \nabla_{\mathbf{r}} \cdot (u_i \nabla_{\mathbf{r}} G - G \nabla_{\mathbf{r}} u_i) - \int_{\mathbf{r} \in \Omega} G \nabla^2 u_i =$$

$$= -\oint_{\mathbf{r} \in \partial \Omega} \hat{\mathbf{n}} \cdot (u_i \nabla_{\mathbf{r}} G - G \nabla_{\mathbf{r}} u_i) + \int_{\mathbf{r} \in \Omega} G(\mathbf{r} - \mathbf{r}_0) f_i(\mathbf{r}, t).$$

An integro-differential boundary problem can be written using boundary conditions. As an example, using Dirichlet and Neumann boundary conditions, the integro-differential problem reads

$$E(\mathbf{r}_0, t) \mathbf{u}(\mathbf{r}_0, t) + \int_{\mathbf{r} \in S_N} \mathbf{u}(\mathbf{r}, t) \, \hat{\mathbf{n}} \cdot \nabla_{\mathbf{r}} G(\mathbf{r} - \mathbf{r}_0) - \int_{\mathbf{r} \in S_D} G(\mathbf{r} - \mathbf{r}_0) \, \hat{\mathbf{n}} \cdot \nabla_{\mathbf{r}} \mathbf{u}(\mathbf{r}, t) =$$

$$= -\int_{\mathbf{r} \in S_D} \mathbf{g}(\mathbf{r}, t) \, \hat{\mathbf{n}} \cdot \nabla_{\mathbf{r}} G(\mathbf{r} - \mathbf{r}_0) + \int_{\mathbf{r} \in S_N} G(\mathbf{r} - \mathbf{r}_0) \, \mathbf{h}(\mathbf{r}, t) + \int_{\mathbf{r} \in \Omega} G(\mathbf{r} - \mathbf{r}_0) \, \mathbf{f}(\mathbf{r}, t).$$

Green's function of the Poisson-Laplace equation reads

$$G(\mathbf{r}; \mathbf{r}_0) = \frac{1}{4\pi} \frac{1}{|\mathbf{r} - \mathbf{r}_0|} \ .$$

### Green's function of the Laplace equation

$$-\nabla^2 G = 0 \qquad \text{for } \mathbf{r} \neq \mathbf{r}_0$$

Solutions with spherical symmetry,

$$0 = \nabla^2 G = \frac{1}{r^2}\left(r^2 G'\right)' \quad \rightarrow \quad G'(r) = \frac{A}{r^2} \quad \rightarrow \quad G(r) = -\frac{A}{r} + B$$

Choosing $B = 0$ s.t. $G(r) \to 0$ as $r \to \infty$, and integrating over a sphere centered in $r = 0$ to get $A = -\frac{1}{4\pi}$,

$$1 = \int_V \delta(r) = -\int_V \nabla^2 G = -\oint_{\partial V} \hat{\mathbf{n}} \cdot \nabla G = -\oint_{\partial V} \hat{\mathbf{r}} \cdot \hat{\mathbf{r}} \frac{A}{r^2} = -4\pi A$$

## 36.2 Helmholtz equation

**todo** from Fourier to Laplace trasnform in the first lines of this section

A Helmholtz's equation can be thouoght as the time Fourier transform of a wave equation,

$$\begin{cases} \dfrac{1}{c^2}\partial_{tt}\mathbf{u}(\mathbf{r}, t) - \nabla^2 \mathbf{u}(\mathbf{r}, t) = \mathbf{f}(\mathbf{r}, t) \\ + \text{ b.c.} \\ + \text{ i.c. }, \end{cases}$$

Fourier transform in time of field $\mathbf{u}(\mathbf{r}, t)$ reads

$$\tilde{\mathbf{u}}(\mathbf{r}, \omega) = \mathcal{F}\{\mathbf{u}(\mathbf{r}, t)\} = \int_{t=-\infty}^{+\infty} \mathbf{u}(\mathbf{r}, t)e^{-i\omega t}\, d\omega$$

and, if $\mathbf{u}(\mathbf{r}, t)$ is compact in time, Fourier transform of its time partial derivatives read

$$\mathcal{F}\{\dot{\mathbf{u}}(\mathbf{r}, t)\} = \int_{t=-\infty}^{+\infty} \dot{\mathbf{u}}(\mathbf{r}, t)e^{-i\omega t}\, d\omega =$$

$$= \mathbf{u}(\mathbf{r}, t)e^{-i\omega t}\Big|_{t=-\infty}^{+\infty} + i\omega \int_{t=-\infty}^{+\infty} \mathbf{u}(\mathbf{r}, t)e^{-i\omega t}\, d\omega =$$

$$= i\omega \mathcal{F}\{\mathbf{u}(\mathbf{r}, t)\}$$

$$\mathcal{F}\{\partial_t^n \mathbf{u}(\mathbf{r}, t)\} = (i\omega)^n \tilde{\mathbf{u}}\,.$$

The differential problem in the transformed domain thus reads

$$-\frac{\omega^2}{c^2}\tilde{\mathbf{u}} - \nabla^2 \tilde{\mathbf{u}} = \tilde{\mathbf{f}}$$

Green's function of Helmholtz'e equation reads

$$G(\mathbf{r}, s) = \alpha^+ \frac{e^{\frac{s|\mathbf{r}-\mathbf{r}_0|}{c}}}{|\mathbf{r} - \mathbf{r}_0|} + \alpha^- \frac{e^{-\frac{s|\mathbf{r}-\mathbf{r}_0|}{c}}}{|\mathbf{r} - \mathbf{r}_0|}$$

with $\alpha^+ + \alpha^- = \frac{1}{4\pi}$.

Being the Laplace transform,

$$\mathcal{L}\{f(t)\} = \int_{t=0^-}^{+\infty} f(t)e^{-st}dt\,,$$

the Laplace transform of a causal function with time delay $\tau \geq 0$ reads

$$\mathcal{L}\{f(t-\tau)\} = \int_{t=0^-}^{+\infty} f(t-\tau)e^{-st}dt = \int_{z=-\tau}^{+\infty} f(z)e^{-s(z+\tau)}\,dz = e^{-s\tau}\int_{z=0}^{+\infty} f(z)e^{-sz}\,dz = e^{-s\tau}\,\mathcal{L}\{f(t)\}$$

having used causality $f(t) = 0$ for $t < 0$. Laplace transform of Dirac's delta $\delta(t)$ reads

$$\mathcal{L}\{\delta(t)\} = \int_{t=0^-}^{+\infty} \delta(t)\,dt = 1\;,$$

so that $e^{-s\tau} = e^{-s\tau}\,1 = \mathcal{L}\{\delta(t-\tau)\}$.

Thus, Green's function for the wave equation reads

$$G(\mathbf{r},t;\mathbf{r}_0,t_0) = \alpha^+ \frac{\delta\left(t-t_0+\frac{|\mathbf{r}-\mathbf{r}_0|}{c}\right)}{|\mathbf{r}-\mathbf{r}_0|} + \alpha^- \frac{\delta\left(t-t_0-\frac{|\mathbf{r}-\mathbf{r}_0|}{c}\right)}{|\mathbf{r}-\mathbf{r}_0|}$$

If $t \geq t_0$, and $G(\mathbf{r},t;\mathbf{r}_0,t_0)$ connects the past $t_0$ with the future $t$, the first term is not causal, and thus $\alpha^+ = 0$ and

$$G(\mathbf{r},t;\mathbf{r}_0,t_0) = \frac{1}{4\pi}\frac{\delta\left(t-t_0-\frac{|\mathbf{r}-\mathbf{r}_0|}{c}\right)}{|\mathbf{r}-\mathbf{r}_0|}\;.$$

## Green's function of Helmholtz's equation

$$\frac{s^2}{c^2}G - \nabla^2 G = \delta(r)$$

$$G(r) = \frac{\alpha e^{kr} + \beta e^{-kr}}{r}$$

Proof:

- Gradient

$$\nabla G(r) = \hat{\mathbf{r}}\partial_r G = \hat{\mathbf{r}}\frac{\alpha(kr-1)e^{kr} + \beta(-kr-1)e^{-kr}}{r^2}$$

- Laplacian

$$\nabla^2 G(r) = \frac{1}{r^2}\left(r^2 G'(r)\right)' =$$

$$= \frac{1}{r^2}\left(\alpha(kr-1)e^{kr} + \beta(-kr-1)e^{-kr}\right)' =$$

$$= \frac{1}{r^2}\left(\alpha k e^{kr} + \alpha k^2 r e^{kr} - \alpha k e^{kr} - \beta k e^{-kr} + \beta k^2 r e^{-kr} + \beta k e^{-kr}\right) =$$

$$= \frac{1}{r}\left(\alpha e^{kr} + \beta e^{-kr}\right)k^2 = k^2 G(r)\;.$$

and thus $k^2 G(r) - \nabla^2 G = 0$, for $r \neq 0$;

- Unity

$$1 = \int_V \delta(r) = \int_V \left(k^2 G - \nabla^2 G\right) = \int_V k^2 G - \oint_{\partial V}\hat{\mathbf{n}}\cdot\nabla G$$

the second term is the sum of two contributions of the form

$$\oint_{\partial V}\hat{\mathbf{n}}\cdot\nabla G^{\pm} = \oint_{\partial V}\frac{\alpha^{\pm}(\pm kr-1)e^{\pm kr}}{r^2} = 4\pi\alpha^{\pm}(\pm kr-1)e^{\pm kr}$$

the first term is the sum of two contributions of the form

$$k^2 \int_V G(r) = k^2 \int_V \frac{\alpha^{\pm} e^{\pm kr}}{r} =$$

$$= k^2 \alpha^{\pm} \int_{R=0}^{r} \int_{\phi=0}^{\pi} \int_{\theta=0}^{2\pi} \frac{e^{\pm kR}}{R} R^2 \sin\phi \, dR \, d\phi \, d\theta =$$

$$= k^2 \alpha^{\pm} \, 4\pi \int_{R=0}^{r} R \, e^{\pm kR} \, dR \, .$$

the last integral can be evaluated with integration by parts

$$\int_{R=0}^{r} R \, e^{\pm kR} \, dR = \left[ \frac{1}{\pm k} e^{\pm kR} R \right]\Big|_{R=0}^{r} \mp \frac{1}{k} \int_{R=0}^{r} e^{\pm kR} \, dR =$$

$$= \frac{1}{\pm k} e^{\pm kr} r - \frac{1}{k^2} e^{\pm kR} + \frac{1}{k^2} =$$

Thus summing everything together,

$$1 = \alpha^{+} \left[ 4\pi k^2 \left( \frac{r}{k} e^{kr} - \frac{1}{k^2} e^{kr} + \frac{1}{k^2} \right) - 4\pi \left( kr - 1 \right) e^{kr} \right] + \alpha^{-} \left[ ... \right] =$$

$$= 4\pi \left( \alpha^{+} + \alpha^{-} \right) \, .$$

## 36.3 Wave equation

Wave equation general problem

$$\begin{cases} \frac{1}{c^2} \partial_{tt} \mathbf{u}(\mathbf{r}, t) - \nabla^2 \mathbf{u}(\mathbf{r}, t) = \mathbf{f}(\mathbf{r}, t) \\ + \text{ b.c.} \\ + \text{ i.c.} \end{cases}$$

Green's problem of the wave equation

$$\frac{1}{c^2} \partial_{tt} G(\mathbf{r}, t; \mathbf{r}_0, t_0) - \nabla_{\mathbf{r}}^2 G(\mathbf{r}, t; \mathbf{r}_0, t_0) = \delta(\mathbf{r} - \mathbf{r}_0)\delta(t - t_0)$$

Integration by parts

$$E(\mathbf{r}_\alpha, t_\alpha) \mathbf{u}(\mathbf{r}_\alpha, t_\alpha) = \int_{t \in T} \int_{\mathbf{r} \in V} \delta(t - t_\alpha)\delta(\mathbf{r} - \mathbf{r}_\alpha) \mathbf{u}(\mathbf{r}, t) =$$

$$= \int_{t \in T} \int_{\mathbf{r} \in V} \left\{ \frac{1}{c^2} \partial_{tt} G - \nabla_{\mathbf{r}}^2 G \right\} \mathbf{u} =$$

$$= \int_{t \in T} \int_{\mathbf{r} \in V} \left\{ \frac{1}{c^2} \left[ \partial_t \left( \mathbf{u} \partial_t G - G \partial_t \mathbf{u} \right) + G \partial_{tt} \mathbf{u} \right] - \nabla_{\mathbf{r}} \cdot \left( \nabla_{\mathbf{r}} G \, \mathbf{u} - G \nabla_{\mathbf{r}} \mathbf{u} \right) - G \nabla_{\mathbf{r}}^2 \mathbf{u} \right\} =$$

$$= \int_{\mathbf{r} \in V} \frac{1}{c^2} \left[ \mathbf{u}(\mathbf{r}, t) \partial_t G(\mathbf{r}, t; \mathbf{r}_\alpha, t_\alpha) - G(\mathbf{r}, t; \mathbf{r}_\alpha, t_\alpha) \partial_t \mathbf{u}(\mathbf{r}, t) \right] \Big|_{t_0}^{t_1} +$$

$$+ \int_{t \in T} \oint_{\mathbf{r} \in \partial V} \left\{ -\hat{\mathbf{n}}(\mathbf{r}, t) \cdot \nabla_{\mathbf{r}} G(\mathbf{r}, t; \mathbf{r}_\alpha, t_\alpha) \, \mathbf{u}(\mathbf{r}, t) + G(\mathbf{r}, t; \mathbf{r}_\alpha, t_\alpha) \hat{\mathbf{n}}(\mathbf{r}, t) \cdot \nabla_{\mathbf{r}} \mathbf{u}(\mathbf{r}, t) \right\} +$$

$$+ \int_{t \in T} \int_{\mathbf{r} \in V} G(\mathbf{r}, t; \mathbf{r}_\alpha, t_\alpha) \underbrace{\left\{ \frac{1}{c^2} \partial_{tt} \mathbf{u}(\mathbf{r}, t) - \nabla_{\mathbf{r}}^2 \mathbf{u}(\mathbf{r}, t) \right\}}_{=\mathbf{f}(\mathbf{r},t)}$$

$$\int_{t \in T} \int_{\mathbf{r} \in V} \frac{1}{4\pi} \frac{\delta\left( t - t_\alpha + \frac{|\mathbf{r} - \mathbf{r}_\alpha|}{c} \right)}{|\mathbf{r} - \mathbf{r}_\alpha|} \mathbf{f}(\mathbf{r}, t) = \int_{\mathbf{r} \in V \cap B_{|\mathbf{r}-\mathbf{r}_\alpha| \leq c(t_\alpha - t)}} \frac{1}{4\pi |\mathbf{r} - \mathbf{r}_\alpha|} \mathbf{f}\left( \mathbf{r}, t_\alpha - \frac{|\mathbf{r} - \mathbf{r}_\alpha|}{c} \right)$$

# Part XII

# Optimization

# OPTIMIZATION

- **Gradient-based methods**, and Hessian based, for continuos functions

- Free and constrained optimization

- …

## 37.1 Unconstrained optimization

Given $f(\mathbf{x})$, with $\mathbf{x} \in \mathbb{R}^n$, a **candidate local** extreme $\mathbf{x}^*$ is a point where the *gradient* $\partial_{\mathbf{x}} f$ of function $f$ is zero

$$\partial_{\mathbf{x}} f(\mathbf{x}^*) = \mathbf{0} .$$

Among these candidates, local extreme may be:

- maximum: if the Hessian evalauted in $\mathbf{x}^*$ is definite negative,

$$\begin{cases} \partial_{\mathbf{x}} f(\mathbf{x}^*) = \mathbf{0} \\ \partial_{\mathbf{xx}} f(\mathbf{x}^*) > 0 \end{cases}$$

- minimum: if the Hessian evalauted in $\mathbf{x}^*$ is definite positive,

$$\begin{cases} \partial_{\mathbf{x}} f(\mathbf{x}^*) = \mathbf{0} \\ \partial_{\mathbf{xx}} f(\mathbf{x}^*) < 0 \end{cases}$$

Here the formalism $\partial_{\mathbf{xx}} f(\mathbf{x}^*) < 0$ indicates that the Hessian matrix is negative definite, see *Example 37.1.1*.

---

**Example 37.1.1 (Positive/negative definite Hessian)**

Local polynomial expansion of a differentiable function $\mathbf{x}$ around $\mathbf{x}^*$ reads

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \Delta\mathbf{x}^T \partial_{\mathbf{x}} f(\mathbf{x}^*) + \frac{1}{2} \Delta\mathbf{x}^T \partial_{\mathbf{xx}} f(\mathbf{x}^*) \, \Delta\mathbf{x} + o(|\Delta\mathbf{x}|^2) ,$$

with $\Delta\mathbf{x} = \mathbf{x} - \mathbf{x}^*$.

Candidate extreme points are thos points $\mathbf{x}^*$ where the gradient vanishes $\partial_{\mathbf{x}} f\mathbf{x}^*) = 0$. If the Hessian $H(\mathbf{x}^*) := \partial_{\mathbf{xx}} f(\mathbf{x}^*)$ is positive definite, by definition,

$$\mathbf{v}^T \, H(\mathbf{x}^*)\mathbf{v} > 0 , \qquad \forall \mathbf{v} \neq \mathbf{0} ,$$

then

$$f(\mathbf{x}) = f(\mathbf{x}^*) + \frac{1}{2} \Delta\mathbf{x}^T H(\mathbf{x}^*) \, \Delta\mathbf{x} + o(|\Delta\mathbf{x}|^2) > f(\mathbf{x}^*) ,$$

as $|\Delta\mathbf{x}| \to 0$, and thus $\mathbf{x}^*$ is a point of minimum, as $f(\mathbf{x}) > f(\mathbf{x}^*)$ for all the neighoring points $\mathbf{x}$.

---

## 37.2 Constrained optimization

$$f(\mathbf{x}) \, , \qquad \mathbf{x} \in D \subseteq \mathbb{R}^n$$

with constraints, such as

$$\mathbf{g}(\mathbf{x}) = \mathbf{0}$$
$$\mathbf{h}(\mathbf{x}) \geq \mathbf{0}$$

### 37.2.1 Method of Lagrange multipliers for equality constraints

$$\tilde{f}(\mathbf{x}; \lambda) = f(\mathbf{x}) - \lambda^T \mathbf{g}(\mathbf{x}) \, .$$

# Part XIII

# Control

# INTRODUCTION TO CONTROL METHODS

**Domain.**

- Frequency domain for linear systems: root locus, Nyquist and Bode criteria,…

- Time domain: *optimal control*, robust optimal control,…

**Full-state, partial-state feedback.**

- Full-state feedback: the complete state of the system can be observed and used for the feedback

- Partial-state feedback

**Observers**, e.g. Kalman filter, tries to estimate the state of a system from its output. Later, this estimatio can be used for feedback.

## 38.1 Optimal control

Optimal control can be recast as a **constrained optimization problem**, $J$, where an extreme - optimum - of an objective function must be found, subject to constraints that include the equations of motion. Some constraints may be included into an augmented objective function $\widetilde{J}$ with the methods of *Lagrange multipliers*.

**Finite time vs. Infinite time horizon.**

### 38.1.1 Generic ODE

$$\mathbf{M}\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u})$$

The objective function combines (weights) the error on a desired performance and the control input, in order to get the desired behavior with feasible control (that can be provided by actuators, without saturation, avoiding unnecessary high power input and too sharp behavior,…)

As an example, if the goal of the control $\mathbf{u}$ is to keep the system around $\mathbf{x} = \mathbf{0}$, the cost function to be minimized can be designed as

$$J = \int_{t=0}^{T} \frac{1}{2} \begin{bmatrix} \mathbf{x}^T & \mathbf{u}^T \end{bmatrix} \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^T & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix} dt \ .$$

## 38.1.2 LTI

$$J = \int_{t=0}^{T} \frac{1}{2} [\mathbf{x}^T \quad \mathbf{u}^T] \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^T & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix} dt$$

$$\begin{cases} \dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu} \\ \mathbf{y} = \mathbf{Cx} + \mathbf{Du} \end{cases}$$

### Infinite-horizon full-state feedback

No need for an observer. The system is assumed to be stable. The augmented cost function reads

$$\widetilde{J}(\mathbf{x}, \mathbf{u}; \lambda) = \int_{t=0}^{+\infty} \left\{ \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{x}^T \mathbf{S} \mathbf{u} + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} - \lambda^T (\dot{\mathbf{x}} - \mathbf{Ax} - \mathbf{Bu}) \right\} dt ,$$

with given initial conditions $\mathbf{x}(0) = \mathbf{x}_0$, so that $\delta \mathbf{x}_0 = \mathbf{0}$.

Using *calculus of variations*, the variations of the cost function w.r.t. $\mathbf{x}$, $\mathbf{u}$, $\lambda$ read

$$\begin{aligned} \delta_{\mathbf{x}} : & \quad \mathbf{Q}\mathbf{x} + \mathbf{S}\mathbf{u} + \dot{\lambda} + \mathbf{A}^T \lambda = \mathbf{0} \\ \delta_{\mathbf{u}} : & \quad \mathbf{S}^T \mathbf{x} + \mathbf{R}\mathbf{u} + \mathbf{B}^T \lambda = \mathbf{0} \\ \delta_{\lambda} : & \quad \dot{\mathbf{x}} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{u} = \mathbf{0} \end{aligned}$$

From the variation w.r.t. $\mathbf{u}$, since $\mathbf{R} > 0$ and thus innvertible,

$$\mathbf{u} = -\mathbf{R}^{-1} \left( \mathbf{S}^T \mathbf{x} + \mathbf{B}^T \lambda \right)$$

Now, assuming the relation $\lambda = \mathbf{P}\mathbf{x}$, it follows

$$\begin{aligned} \dot{\lambda} &= -\mathbf{Q}\mathbf{x} - \mathbf{S}\mathbf{u} - \mathbf{A}^T \lambda = \\ &= \left\{ -\mathbf{Q} + \mathbf{S}\mathbf{R}^{-1} \left( \mathbf{S}^T + \mathbf{B}^T \mathbf{P} \right) - \mathbf{A}^T \mathbf{P} \right\} \mathbf{x} \\ \dot{\lambda} &= \dot{\mathbf{P}}\mathbf{x} + \mathbf{P}\dot{\mathbf{x}} = \\ &= \dot{\mathbf{P}}\mathbf{x} + \mathbf{P} \left( \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \right) = \\ &= \left\{ \dot{\mathbf{P}} + \mathbf{P} \left[ \mathbf{A} - \mathbf{B}\mathbf{R}^{-1} \left( \mathbf{S}^T + \mathbf{B}^T \mathbf{P} \right) \right] \right\} \mathbf{x} , \end{aligned}$$

and comparing the two different expressions of $\dot{\lambda}$, if the equality holds for any $\mathbf{x}$, the **dynamical Riccati equation** for $\mathbf{P}$ is derived as

$$\dot{\mathbf{P}} + \mathbf{P}\widetilde{\mathbf{A}} + \widetilde{\mathbf{A}}^T \mathbf{P} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \mathbf{P} + \widetilde{\mathbf{Q}} = \mathbf{0} ,$$

where $\widetilde{\mathbf{A}} = \mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{S}^T$ and $\widetilde{\mathbf{Q}} = \mathbf{Q} - \mathbf{S}\mathbf{R}^{-1}\mathbf{S}^T$. Riccati equation is a non-linear dynamical matrix equation in $\mathbf{P}$. Algorithms for computing the solution of dynamical and algebraic equation exists, see *Example 38.1.1*.

Once $\mathbf{P}$ is evaluated, the control law reads

$$\mathbf{u} = -\mathbf{R}^{-1} \left( \mathbf{S}^T + \mathbf{B}^T \mathbf{P} \right) \mathbf{x} .$$

For infinite-horizone, the algebraic equation (ARE) for the steady state is solved after setting $\dot{\mathbf{P}} = \mathbf{0}$, the solution for a LTI system is a constant matrix $\mathbf{P}$, and thus the control law is a proportional feedback on the full-state of the system,

$$\mathbf{u} = -\mathbf{G}\mathbf{x} ,$$

with $\mathbf{G} = \mathbf{R}^{-1} \left( \mathbf{S}^T + \mathbf{B}^T \mathbf{P} \right)$.

**Example 38.1.1 (Solution of Riccati equation)**

…

**Properties. todo**

- **P** symmetric? definite positive? …

- …

# Part XIV

# Reinforcement Learning

# INTRODUCTION TO REINFORCEMENT LEARNING

# MARKOV PROCESSES

---

**Definition 40.1 (Stochastic process)**

A stochastic process is as collection of random variables $X(t)$ defined on a common probability space $(\Omega, \mathcal{F}, P)$, indexed by some set $T$, with outcomes in a measurable space $(S, \Sigma)$,

$$\{X(t) : t \in T\} \ .$$

---

Markov process is a memomry-less stochastic process (see def *Definition 40.1*). For discrete-time Markov processes, state at time $n$ only depends on state at time $n - 1$

$$\langle S, P, \mu \rangle$$

with

- $S$ (finite) set of states

- **P** state transition probability matrix, $P_{ss'} = P(s'|s)$

- $\mu$ a set of initial probability, $\mu_i = P(X_0 = i), \forall i$

**Properties of probability matrix.**

A **stationary MP** is a process described by a constant state transition probability.

**1-step transition.**

---

**Notation: probability vectors as row vectors**

In treating stochastic processes, probability vectors are usually treated as row vectors. Probability distribution over states $s'$ at time $n$ is the sum of probabilities or reaching $s'$ starting from $s$ at time $n - 1$

$$p_n(s') = \sum_s P(s'|s)p_{n-1}(s) = \sum_s p_{n-1}(s)P_{ss'} \ ,$$

or using matrix formalism

$$\mathbf{p}_n = \mathbf{p}_{n-1}\mathbf{P} \ .$$

---

**Properties of probability matrix.**

Since $P_{ss'} = P(s'|s)$ represents the probability of arriving in $s'$ starting from $s$ (i.e. the *conditional probability* of arriving in $s'$, starting from $s$), the entries of a matrix representing a time-discrete process are

$$0 \leq P_{ss'} \leq 1 \ .$$

---

The sum over $s'$ is 1 (sum over all the possibilities),

$$\sum_{s'} P(s'|s) = \sum_{s'} P_{ss'} = 1 , \quad \forall s \quad \text{or} \quad \mathbf{P1} = \mathbf{1} .$$

Thus, a probability matrix $\mathbf{P}$ has a eigenvalue $\lambda = 1$ with the uniform vector $\mathbf{1}$ as eigenvector.

**Is this unique?** *No, in general.*

**Existence of a stationary distribution**

*Spectral radius* of a probability matrix is $\rho(\mathbf{P}) = 1$.

**Proof**

By *Gershgorin circle theorem*,

$$|\lambda - P_{ii}| \le \sum_{j \ne i} |P_{ij}| ,$$

for $i$ s.t. the components of the (right) eigenvalue $\mathbf{v}$ satisfy $|x_i| \ge |x_j|, \forall j$. As $P_{ij} \ge 0$ and $\sum_j P_{ij} = P_{ii} + \sum_{j \ne i} P_{ij} = 1$, it follows

$$|\lambda - P_{ii}| \le \sum_{j \ne i} P_{ij} = \sum_j P_{ij} - P_{ii} = 1 - P_{ii} ,$$

meaning that any eigenvalue $\lambda$ is contained in a circle centered in $P_{ii}$ with radius $1 - P_{ii}$. It's not hard to prove that such a circle is contained in a circle centered in $0$ with radius $1$[1]. Thus, for all the eigenvalues $|\lambda| \le 1$ holds, and at least one eigenvalue with $|\lambda| = 1$, namely the $\lambda = 1$ eigenvalue shown above, exists. Thus, the spetral radius of a probability matrix is $\rho(\mathbf{P}) = 1$.

**n-step transition.** $n$-step transition in a stationary MP reads

$$\mathbf{p}_{m+n} = \mathbf{p}_{m+n-1}\mathbf{P} = \mathbf{p}_{m+n-2}\mathbf{PP} = \cdots = \mathbf{p}_m \mathbf{P}^n .$$

**First passage.** …; recurrent states

**Classification of states**: accessible, communicating (and dividing states into partitions, disjoint classes); positive recurrent, periodic, ergodic; absorbing or transient

**Classification of MP**: absorbing, ergodic, regular

**Stationary distribution,** a probability vector $\mathbf{p}$ so that probability doesn't change in the next timestep,

$$\mathbf{p} = \mathbf{p}\mathbf{P} .$$

**Fundamental matrix**

**Other definitions**: mixing rate, spectral gap

# 40.1 Markov Reward Processes

$$\langle S, P, R, \gamma, \mu \rangle$$

with

---

[1] $1 - P_{ii} \ge |\lambda - P_{ii}| = |\lambda - 0 + 0 - P_i i| \ge ||\lambda| - P_i i|$; if $|\lambda| \ge P_{ii}$ it follows $1 \ge |\lambda|$; if $|\lambda| \le P_{ii}$…**todo**

- $R$ **reward function**; it could be a stochastic variable, depending on a state with probability $p(r|s)$ *(or on the transition $p(r|s, s')$?)*

- $\gamma$ **discount factor**, $\gamma \in [0, 1]$; it's used to represent the present value of future rewards ($\gamma \sim 0$ short-sighted, $\gamma \sim 1$ far-sighted)

**Return**

- Time horizion: finite, indefinite (until stopping criteria, or absorbing state), infinite

- Possible choices of cumulative reward:

  - total reward, $V = \sum_i r_i$

  - average reward, $V = \frac{1}{n} \sum_{i=1}^{n} r_i$

  - discounted reward $V = \sum_i \gamma^{i-1} r_i$

Return $v_t$ is defined as the total discount reward from timestep $t$

$$v_t = \sum_{k=0}^{+\infty} \gamma^k r_{t+k+1} = r_{t+1} + \gamma r_{t+2} + ...$$

**Value function** Expected value of the return $v_t$ from state $s_t = s$

$$V(s) := \mathbb{E}\left[v_t | s_t = s\right]$$

### 40.1.1 Bellman Equation for $V(s)$

$$
\begin{aligned}
V(s) &= \mathbb{E}\left[v_t | s_t = s\right] = \\
&= \mathbb{E}\left[ r_{t+1} + \gamma \sum_{k=0}^{+\infty} \gamma^k r_{t+2+k} \middle| s_t = s \right] = \\
&= \mathbb{E}\left[ r_{t+1} + \gamma v_{t+1} | s_t = s \right] = \\
&= \mathbb{E}\left[ r_{t+1} | s_t = s \right] + \gamma \mathbb{E}\left[ v_{t+1} | s_t = s \right] = \\
&= R(s) + \gamma \sum_{s'} P(s'|s) \mathbb{E}\left[ v_{t+1} | s_{t+1} = s' \right] = \\
&= R(s) + \gamma \sum_{s'} P(s'|s) V(s') \,.
\end{aligned}
$$

or in matrix form

$$\mathbf{V} = \mathbf{R} + \gamma \mathbf{P} \mathbf{V} \,.$$

**Solution of Bellman equation.** Bellman equation for value function is linear. It can be solved directly, solving the linear problem

$$(\mathbf{I} - \gamma \mathbf{P}) \mathbf{V} = \mathbf{R} \,.$$

Solutions:

- direct solution of the linear system (feasible for small-dimensional MRP only?)

- iterative methods for large MRPs: DP (dynamic programming), MC (Monte-Carlo evaluation), TD (Temporal-Difference learning)

## 40.2 Markov Decision Processes

$$\langle S, A, P, R, \gamma, \mu \rangle$$

with

- A (finite) set of actions
- **P** state transition probability matrix, $P(s'|s, a)$
- $R$ reward function, $R(s, a) = \mathbb{E}\left[r|s, a\right]$

### 40.2.1 Policies

A policy is a distribution over actions, given the initial state

$$\pi(a|s) = P(a|s)$$

Given a policy $\pi(a|s)$, a MDP becomes the MRP with

- transition probability

$$P^\pi(s'|s) = \sum_{a \in A} \pi(a|s) P(s'|s, a)$$

- reward (**todo** *Prove it!*)

$$R^\pi(s) = \sum_{a \in A} \pi(a|s) R(s, a)$$

**Value functions.**

- State value function, $V^\pi(s)$, expected return starting from $s$ and following policy $\pi$

$$V^\pi(s) := \mathbb{E}_\pi\left[v_t|s_t = s\right]$$

- Action-value functon $Q^\pi(s, a)$, expected return starting from $s$, taking action $a$ first and then following policy $\pi$

$$Q^\pi(s, a) := \mathbb{E}_\pi\left[v_t|s_t = s, a_t = a\right]$$

## 40.2.2 Bellman Equations

**Bellman expectation equations.** Bellman expectation equation for the state-value function $V^\pi(s)$ reads

$$
\begin{aligned}
V^\pi(s) &:= \mathbb{E}_\pi\left[v_t | s_t = s\right] = \\
&= \mathbb{E}_\pi\left[ r_{t+1} + \gamma \sum_{k=0}^{+\infty} \gamma^k r_{t+2+k} \bigg| s_t = s \right] = \\
&= \sum_{a \in A} \pi(a|s) \mathbb{E}_\pi\left[ r_{t+1} + \gamma \sum_{k=0}^{+\infty} \gamma^k r_{t+2+k} \bigg| s_t = s, a_t = a \right] = \\
&= \sum_{a \in A} \pi(a|s) \mathbb{E}\left[ r_{t+1} | s_t = s, a_t = a \right] + \gamma \sum_{a \in A} \sum_{s' \in S} \pi(a|s) P(s'|s, a) \mathbb{E}_\pi\left[ v_{t+1} | s_{t+1} = s' \right] = \\
&= \sum_{a \in A} \pi(a|s) R(s, a) + \gamma \sum_{s' \in S} P^\pi(s'|s) \mathbb{E}_\pi\left[ v_{t+1} | s_{t+1} = s' \right] = \\
&= \sum_{a \in A} R(s, a) \left\{ \pi(a|s) + \gamma \sum_{s' \in S} P(s'|s, a) V^\pi(s') \right\} = \\
&= R^\pi(s) + \gamma \sum_{s' \in S} P^\pi(s'|s) V^\pi(s') =
\end{aligned}
\tag{40.1}
$$

Bellman expectation equation for the action-value function $Q^\pi(s)$ reads

$$
\begin{aligned}
Q^\pi(s, a) &:= \mathbb{E}_\pi\left[v_t | s_t = s, a_t = a\right] = \\
&= \mathbb{E}\left[ r_{t+1} | s_t = s, a_t = a \right] + \gamma \mathbb{E}_\pi\left[ v_{t+1} | s_t = s, a_t = a \right] = \\
&= R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \mathbb{E}_\pi\left[ v_{t+1} | s_{t+1} = s' \right] = \\
&= R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^\pi(s') = \\
&= R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \sum_{a' \in A} \pi(a'|s') \mathbb{E}_\pi[v_{t+1}|s_{t+1} = s', a_{t+1} = a'] = \\
&= R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \sum_{a' \in A} \pi(a'|s') Q^\pi(s', a')
\end{aligned}
\tag{40.2}
$$

**Bellman expectation operators todo**

**Optimal value functions, and optimal policies**

---

**Definition 40.2.1 (Optimal value functions)**

Optimal state-value function $V^*(s)$ is defined as

$$
V^*(s) = \max_\pi V^\pi(s) \,.
$$

Optimal action-value function $Q^*(s, a)$ is defined as

$$
Q^*(s, a) = \max_\pi Q^\pi(s, a) \,.
$$

---

Usually, the goal of a problem involving MDP is the maximization of value functions, as the performance of the MDP.

Valus functions define a partial ordering of policies,

$$
\pi \geq \pi' \text{ if } V^\pi(s) \geq V^{\pi'} \ \forall s \in S \,.
$$

**Theorem 40.2.1 (Optimal policies)**

For a MDP,

- an optimal policy $\pi^*$ exists, s.t. $\pi^* \geq \pi$, $\forall \pi$

- all optimal policies achieve the same value of optimal functions, namely

$$V^{\pi^*}(s) = V^*(s) \quad , \quad Q^{\pi^*}(s,a) = Q^*(s,a) \ .$$

- a deterministic optimal policy exists, and it's found optimizing action-value function over actions $a$,

$$\pi^*(a|s) = \begin{cases} 1 & , \quad \text{if } a = \text{argmax}_{a \in A} Q^*(s,a) \\ 0 & , \quad \text{otherwise} \end{cases}$$

## Proof

**todo**

**Property 40.2.1 (Property: optimal state-value and action-value functions)**

$$V^*(s) = \max_a Q^*(s,a) \tag{40.3}$$

**todo** *Prove it!*

**Bellman optimality equations.** Using *Property 40.2.1* and the expression of action-value function in (40.2)

$$V^*(s) = \max_a Q^*(s,a) =$$

$$= \max_a \max_\pi Q^\pi(s,a) =$$

$$= \max_a \max_\pi \left\{ R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V^\pi(s') \right\} =$$

$$= \max_a \left\{ R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V^*(s') \right\}$$

$$Q^*(s,a) = \max_\pi Q^\pi(s,a) =$$

$$= \max_\pi \left\{ R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V^\pi(s') \right\} =$$

$$= R(s,a) + \max_\pi \left\{ \gamma \sum_{s' \in S} P(s'|s,a) V^\pi(s') \right\} =$$

$$= R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V^*(s') =$$

$$= R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} Q^*(s',a')$$

**Bellman optimality operators**

**Properties of Bellman operators**

**Solving Bellman optimality equations**

- non-linear

- no closed form solution

- iterative solutions:

    – Dynamic Programming (DP): Value Iteration (VI), Policy Iteration (PI)

    – Linear Programming (LP)

    – Reinforcement Learning (RL): Q-learning (off-policy), SARSA (on-policy),…

# FORTYONE

# METHODS OF SOLUTION OF MPD: DP AND LP

## 41.1 Brute force: policy search

## 41.2 Dynamic programming

### Introduction to DP

- sequential or temporal approach to optimization
- solving complex problems breaking them down into sub-problems:
    - solve sub-problems
    - combine solutions of sub-problems
- for problems with 2 properties:
    - optimal substructure: optimal sol can be decomposed in sub-pbs
    - overlapping sub-problems: sub-pbs recur, sols can be cached and reused
- DP assumes full knowledge of the MDP
- used for planning (???)
- prediction: given MDP and policy, evaluate value function
- control: given MDP, evaluate optimal value function and policy

### 41.2.1 Policy iteration

Policy iteration alternates *policy evaluation* and *policy improvement* steps. When value function converges, iteration stops at **Bellman optimality condition**, as shown in the *algorithm* below.

## Policy Evaluation

**Bellman equation for a given policy $\pi$.** Solution of Bellman equation (40.1),

$$V^\pi(s) = \sum_{a \in A} R(s,a) \left\{ \pi(a|s) + \gamma \sum_{s' \in S} P(s'|s,a) V^{pi}(s') \right\}$$

- direct solution. Complexity: …

- full policy-evaluation back-up

$$V_{k+1}(s) \leftarrow \sum_{a \in A} R(s,a) \left\{ \pi(a|s) + \gamma \sum_{s' \in S} P(s'|s,a) V^\pi_k(s') \right\}$$

synchronous, asynchronous update; **todo** *convergence?*

## Policy Improvement

**Greedy update.** Let $\pi$ be a deterministic policy. Greedy update acts maximizing the action-value function for all state over all the actions,

$$\pi'(s) := \mathrm{argmax}_{a \in A} Q^\pi(s,a) \ ,$$

and it "improves the value function" (which value function?) of any state,

$$Q^\pi(s, \pi'(s)) := \max_{a \in A} Q^\pi(s,a) \geq Q^\pi(s, \pi(s)) = V^\pi(s) \ . \tag{41.1}$$

In stochastic optimization, usually *$\varepsilon$-greedy improvement* is introduced to improve **exploration** and try to avoid converging to local optimum.

---

**Theorem 41.2.1 (Policy improvement theorem)**

Let $\pi$ and $\pi'$ a pair of deterministic policites s.t.

$$Q^\pi(s, \pi'(s)) \geq V^\pi(s) \ , \quad \forall s \in S$$

Then $\pi' \geq \pi$, i.e. $V^{\pi'}(s) \geq V^\pi(s), \forall s \in S$.

---

**Proof.**

$$
\begin{aligned}
V^\pi(s) \leq Q^\pi(s, \pi'(s)) &= \\
&=: \mathbb{E}_\pi \left[ v_t | \, s_t = s, a_t = \pi'(s) \right] = \\
&= \mathbb{E}_\pi \left[ r_{t+1} + \gamma v_{t+1} | \, s_t = s, a_t = \pi'(s) \right] = \\
(1) &= \mathbb{E}_{\pi'} \left[ r_{t+1} | s_t = s \right] + \gamma \mathbb{E}_\pi \left[ v_{t+1} | \, s_{t+1} = s', a_{t+1} = \pi(s) \right] = \\
&= \mathbb{E}_{\pi'} \left[ r_{t+1} | s_t = s \right] + \gamma V^\pi(s) = \\
&\leq \mathbb{E}_{\pi'} \left[ r_{t+1} | s_t = s \right] + \gamma Q^\pi(s, \pi'(s)) = \\
&= \mathbb{E}_{\pi'} \left[ r_{t+1} | s_t = s \right] + \gamma \mathbb{E}_\pi \left[ r_{t+2} + \gamma v_{t+2} | \, s_{t+1} = s, a_{t+1} = \pi'(s) \right] = \\
(2) &= \mathbb{E}_{\pi'} \left[ r_{t+1} | s_t = s \right] + \gamma \mathbb{E}_{\pi'} \left[ r_{t+2} | s_{t+1} = s \right] + \gamma^2 \mathbb{E}_\pi \left[ v_{t+2} | \, s_{t+2} = s'', a_{t+2} = \pi(s) \right] = \\
&= \mathbb{E}_{\pi'} \left[ r_{t+1} + \gamma r_{t+2} | s_t = s \right] + \gamma^2 V^\pi(s) = \\
&\leq \mathbb{E}_{\pi'} \left[ r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + ... | s_t = s \right] = \\
&\leq V^{\pi'}(s) \ .
\end{aligned}
$$

---

with in (1) the first action $a_t$ is drawn from $\pi'(s)$, $s'$ results from the transition

$$P^\pi(s'|s) = P(s'|s,a)\pi(a|s) \ ,$$

and the following are drawn from $\pi$. The same notation is used in (2), and in all the following manipulations.

### Policy Iteration algorithm

Policy iteration algorithm

$$\pi_0 \to V^{\pi_0} \to \pi_1 \to V^{\pi_1} \to ... \pi^* \to V^* \to \pi^*$$

Convergence is reached when $\pi_{n+1}(s) = \pi_n(s)$. Greedy update (41.1) gives

$$Q^{\pi_n}(s, \pi_{n+1}(s)) = Q^{\pi_n}(s, \pi_n(s)) = V^{\pi_n}(s) \ ,$$

namely, when the optimality condition (40.3) is reached.

## 41.2.2 Value Iteration

Unlike policy iteration, value iteration algorithm doesn't provide any explicit policy, but only optimal value of the value-functions.

An optimal policy for a MDP is built with an optimal action $a^*$ first, followed by actions drawn from the optimal policy. Thus, it must satisfy

$$V^*(s) = \max_{a\in A} R(s,a) + V^*(s') \ ,$$

with $s_n = s$ and $s_{n+1} = s'$, resulting from transition $s'|s,a$. Subdividing the problem in sub-problems, **deterministic value iteration** proceeds from $s'$ backwards to $s$,

$$V^*(s) \leftarrow \max_{a\in A} R(s,a) + V^*(s') \ .$$

or

$$V^*(s) \leftarrow \max_{a\in A} \left\{ R(s,a) + \gamma \sum_{s'\in S} P(s'|s,a)V(s') \right\} \ .$$

## 41.2.3 Dynamic programming: extensions

- synchronous, or in-place
- prioritized sweeping/update
- ...

# 41.3 Linear programming

...

# FORTYTWO

# METHODS OF SOLUTION OF MPD: RL

## Introduction

- model-free / model-based

- on-policy / off-policy

- online / offline

- tabular / function approximation

- value-based / policy-based / actor-critic

Two main goals of methods in RL:

- *prediction or evaluation*: estimate the performance of a given policy

- *control*: find the best (or a good policy) to get the best performance

These two tasks usually rely on two processes: *evaluation* and *improvement*. Evaluation stage provides the perfomance of a given policy usually in terms of value functions. Improvement aims at finding a new policy with better performance.

Prediction involves only the evaluation step, while control usually involves an iteration over alternate evaluation and improvement processes.

## 42.1 Evaluation

### 42.1.1 Monte-Carlo RL

- model-free, learn from experience of complete episodes ((-) no bootstrapping)

- uses empirical mean return, instead of expected value

- first-visit MC (unbiased estimator) / every-visit MC (biased estimator)

**First-visit MC**

Initialize: $\pi$ policy to be evaluated; $V$ arbitrary state-value function, $R(s)$ empty list of returns

Loop until convergence or stopping criterion

- generate an episode using $\pi$

  - for each state $s$ in the episode:

    * evalaute and append the return $r(s)$ following the first occurrence of $s$ to the list of $R(s)$

    * update $V(s) = \text{average}(R(s))$

**Every-visit MC**

Initialize: $\pi$ policy to be evaluated; $V$ arbitrary state-value function, $R(s)$ empty list of returns

Loop until convergence or stopping criterion

- generate an episode using $\pi$
    - for each state $s$ in the episode:
        * for each occurrence of state $s$ in the episode:
            · evaluate and append the return $r(s)$ following the occurrence of $s$ to the list of $R(s)$
            · update $V(s) = \text{average}(R(s))$

**Using incremental mean update**

$$\begin{aligned}
\mu_{N+1} &= \frac{x_1 + \cdots + x_N + x_{N+1}}{N+1} = \\
&= \frac{x_1 + \cdots + x_N}{N} \frac{N}{N+1} + \frac{x_{N+1}}{N+1} = \\
&= \frac{N}{N+1}\mu_N + \frac{1}{N+1}x_{N+1} = \\
&= \mu_N + \frac{1}{N+1}\left(x_{N+1} - \mu_N\right)
\end{aligned}$$

Incremental update of the average of the value function reads

$$V_k(s) \leftarrow V_{k-1}(s) + \frac{1}{k}\left(v_k - V_{k-1}(s)\right) \ ,$$

with $v_k$ the return of state $s_k$.

A generalized update may follow, to introduce a free hyperparameter as a weight of older estimation,

$$\begin{aligned}
V_k(s) &\leftarrow V_{k-1}(s) + \alpha_k\left(v_k - V_{k-1}(s)\right) \\
&= (1 - \alpha_k)V_{k-1}(s) + \alpha_k v_k
\end{aligned}$$

## 42.1.2 Temporal Difference (TD) Learning

- model-free, learn from experience of incomplete episodes (bootstrapping)
- …

## 42.1.3 TD prediction

- Starting from every visit MC

$$V(s_t) \leftarrow V(s_t) + \alpha(v_t - V(s_t))$$

- Update towards an **estimation of the return**
    - the simplest TD algorithm is $TD(0)$.
        * estimated return (called **TD target**) reads

$$v_t \sim r_{t+1} + \gamma V(s_{t+1})$$

* **TD error** error $\delta_t$ is defined as

$$\delta_t := r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

* update step

$$V(s) \leftarrow V(s) + \alpha_k \left( r_{t+1} + \gamma V(s_{t+1}) - V(s_{t+1}) \right) =$$
$$= V(s) + \alpha_k\, \delta_t$$

– $TD(n)$ algorithm uses $n$-step prediction

$$v_t^{(n)} = r_{t+1} + \gamma r_{t+2} + ... \gamma r_{t+n} + \gamma^{n+1} V(s_{t+n})$$

– $\lambda$-return $v_t^\lambda$ combines all $n$-step returns $v_t^{(n)}$ as

$$v_t^\lambda := (1 - \lambda) \sum_{n=1}^{+\infty} \lambda^{n-1} v_t^{(n)}$$

weights s.t. sum of weights $= 1$ and decay as $\lambda(< 1)$

– …eligibility traces, forward and backward TD,…

## 42.2 Model control

### 42.2.1 On- and Off-Policy Learning

* On-policy: learn policy $\pi$ from experience sampled from $\pi$
* Off-policy: learn policy $\pi$ from experience sampled from $\tilde{\pi}$

**Greedy policy improvement**

* over $V(s)$ requires model of MDP,

$$\pi'(s) = \operatorname{argmax}_{a \in A} \left\{ R(s, a) + P(s'|s, a) V(s') \right\}$$

* **over $Q(s, a)$ is model-free**

$$\pi'(s) = \operatorname{argmax}_{a \in A} Q(s, a) \ .$$

### 42.2.2 MC control

**Note:** Generalized policy iteration with MC evaluation

MC evaluation is model-free. To keep a model-free control task, a model-free policy improvement is required, like greedy policy improvement over $Q(s, a)$.

## 42.2.3 $\epsilon$-greedy policy improvement

$$\pi(a|s) = \begin{cases} 1 - \varepsilon \left(1 - \dfrac{1}{N}\right) & , & \text{if } a = \text{argmax}_{a \in A} Q^\pi(s, a) \\ \dfrac{\varepsilon}{N} & , & \text{otherwise} \end{cases}$$

**Theorem 42.2.1 ($\varepsilon$-greedy Policy Improvement)**

For ..., the $\varepsilon$-greedy policy $\pi'$ w.r.t. $Q^\pi$ is an improvement, $V^{\pi'}(s) \geq V^\pi(s)$.

**Proof**

$$\begin{aligned} Q^\pi(s, \pi'(s)) &= \sum_{a \in A} \pi'(a|s) Q^\pi(s, a) = \\ &= \sum_{a \in A, a \neq a^*} \frac{\varepsilon}{m} Q^\pi(s, a) + \left(1 - \varepsilon + \frac{\varepsilon}{m}\right) Q^\pi(s, a^*) = \\ &= \sum_{a \in A} \frac{\varepsilon}{m} Q^\pi(s, a) + (1 - \varepsilon) Q^\pi(s, a^*) = \\ (1) &\geq \sum_{a \in A} \frac{\varepsilon}{m} Q^\pi(s, a) + (1 - \varepsilon) \sum_{a \in A} \frac{\pi(a|s) - \frac{\varepsilon}{m}}{1 - \varepsilon} Q^\pi(s, a) = \\ &= \sum_{a \in A} \pi(a|s) Q^\pi(s, a) = V^\pi(s) \end{aligned}$$

having used in $(1)$ ... And for the policy improvement theorem *Theorem 41.2.1*, $V^{\pi'}(s) \geq V(s)$. **todo** *Check it!*

**GLIE (Greedy in the Limit of Infinite Exploration)** ...

**GLIE MC Control**

- $k^{th}$ episode using \pi$$

- for each state $s_t$, and action $a_t$ in the episode, update number of $(s, a)$ pair occurence and action-value function

$$\begin{aligned} N(s_t, a_t) &\leftarrow N(s_t, a_t) + 1 \\ Q(s_t, a_t) &\leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)}(v_t - Q(s_t, a_t)) \end{aligned}$$

- improve policy with $\varepsilon$-greedy method over action-value function $Q$,

$$\begin{aligned} \varepsilon &\leftarrow \frac{1}{k} \\ \pi &\leftarrow \varepsilon - \text{greedy}(Q) \end{aligned}$$

### 42.2.4 TD control

- (+): online, from incomplete episodes, lower variance

- use TD instead of MC for policy evaluation, **SARSA**:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left( r + \gamma Q(s',a') - Q(s,a) \right)$$

…

**Variations.** SARSA($\lambda$),…

---

**Note:** SARSA is on-policy

---

### 42.2.5 Off-policy learning

**Q-learning**, is off-policy as $a'$ is not taken from $\pi$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left( r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right)$$

## 42.3 Actor-Critic

**todo**

## 42.4 Learning and planning

- Model-free RL: no model; **learn** policy and/or value function from real experience

- Model-based RL: learn a model from real experience; **plan** policy and/or value function from simulated experience

- Dyna: leanr a model from real experience; **learn and plan** a policy and/or value function from real and simulated experience

# FORTYTHREE

# LARGE OR CONTINUOUS MDPS

For small-dimensional problems, value functions can be efficiently represented by arrays or look-up tables,

$$V(s) \quad \rightarrow \quad \mathbf{V} \in \mathbb{R}^{|S|}$$
$$Q(s,a) \quad \rightarrow \quad \mathbf{Q} \in \mathbb{R}^{|S| \times |A|}$$

For large problems with large number of discrete states $\sim 10^N$ ($N = 20$ for backgammon, $= 10^{170}$ for Go), or continuous state space, RL usually relies on **function approximation** of the value function,

$$V^\pi(s) \sim V(s; \theta)$$
$$Q^\pi(s, a) \sim Q(s, a; \theta)$$

Parameters $\theta$ are updated during learning. Examples of function approximators, the most suitable depend on the problem itself (ANN, Fourier/wavelets basis, polynomial, piecewise-polynomial, coarse coding,…)

**Optimizing/minimizing mean-squared error between** $V(s, \theta)$ **and** $V^\pi(s)$

$$L(\theta) := \mathbb{E}_\pi \left[ \left( V^\pi(s) - V(s; \theta) \right)^2 | s_t = s \right]$$

Gradient descent,

$$\Delta \theta = -\frac{1}{2} \alpha \nabla_\theta L(\theta) =$$
$$= -\frac{1}{2} \alpha \nabla_\theta \, \mathbb{E}_\pi \left[ \left( V^\pi(s) - V(s, \theta) \right)^2 | s_t = s \right] =$$
$$= \alpha \, \mathbb{E}_\pi \left[ \left( V^\pi(s) - V(s, \theta) \right) \nabla_\theta V(s, \theta) \mid s_t = s \right] ,$$

and the stochastic gradient decent update reads

$$\Delta \theta = \alpha \, \left( V^\pi(s) - V(s, \theta) \right) \, \nabla_\theta V(s, \theta) .$$

State representation with **features** $\phi(s)$,

$$\phi(s) = \begin{bmatrix} \phi_1(s) \\ ... \\ \phi_n(s) \end{bmatrix}$$

**State-value function as a linear combination of features**

$$V(s; \theta) = \phi^T(s) \, \theta ,$$

and the gradient of the objective function directly follows from

$$\nabla_\theta V(s; \theta) = \phi(s) .$$

…

**Action-value function as a linear combination of features**

$$Q(s, a; \theta) = \phi^T(s, a) \, \theta ,$$

…

## 43.1 Value and Policy-Based RL

- value based

- policy based

- actor-critic

…