

# Numerical Algorithms

Fabian Bosshard

February 6, 2026

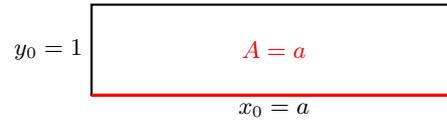
## Contents

<b>1</b>	<b>Introduction</b> .....	<b>1</b>
1.1	Root Finding .....	1
<b>2</b>	<b>Iterative Methods</b> .....	<b>2</b>
2.1	Jacobi Method .....	2
2.2	Gauss–Seidel Method .....	2
<b>3</b>	<b>Polynomial Interpolation</b> .....	<b>3</b>
3.1	Polynomial Curves .....	3
3.2	Neville’s Algorithm .....	4
3.3	Lagrange Form .....	4
3.4	Barycentric Form .....	5
3.5	Newton Form and Divided Differences .....	5
<b>4</b>	<b>Bézier Curves and Splines</b> .....	<b>7</b>
4.1	Bézier Curves .....	7
4.2	Bézier Splines .....	9
4.2.1	$C^0$ -continuity .....	9
4.2.2	$C^1$ -continuity .....	10
4.2.3	$C^2$ continuity .....	10
4.3	B-Splines .....	10
4.3.1	Definition and evaluation .....	11
4.3.2	B-spline basis functions .....	11
4.3.3	Interpolation .....	12
<b>5</b>	<b>Linear Least Squares</b> .....	<b>13</b>
5.1	Data Fitting (Functional Models) .....	13
5.2	Curve Fitting (Parametric Models) .....	14
5.3	QR Factorization .....	15
5.3.1	Classical Gram-Schmidt .....	15
5.3.2	Solving least squares problems .....	16
5.3.3	Modified Gram–Schmidt and Householder .....	16
5.4	Weighted Least Squares .....	16
5.5	Least-Norm Solutions .....	17
5.6	Pseudoinverse .....	18
<b>6</b>	<b>Nonlinear Least Squares</b> .....	<b>19</b>

# 1 Introduction

## 1.1 Root Finding

**Goal:** Find  $\sqrt{a}$  for  $a > 0$ .



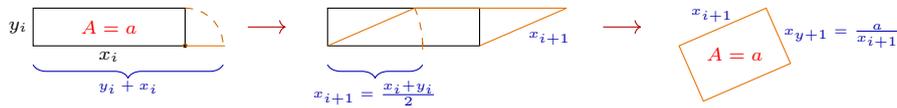
**Idea:** successively transform conserving area until perfect square.

**Straightedge and compass approach:**

$$x_{i+1} = \frac{x_i + y_i}{2}$$

Since  $A = x_i \cdot y_i = a \implies y_i = \frac{a}{x_i}$  (can always express like this):

$$x_{i+1} = \frac{x_i + \frac{a}{x_i}}{2} = \frac{x_i^2 + a}{2x_i} = x_i - \frac{x_i^2 - a}{2x_i}$$



special case of **Newton's Method**, i.e. finding  $x \in \mathbb{R}$  such that  $f(x) = 0$ :

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

**Remark:**  $|x - x_i|$  decreases quadratically in  $i$  (can be shown using Taylor expansion).

**BUT careful:** can also fail badly if

- $x_0$  is too far from correct solution  $x$
- $f'(x_i) \approx 0$  for some  $x_i$
- $f'$  unbounded close to the root  $x$

## 2 Iterative Methods

### 2.1 Jacobi Method

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right), \quad i = 1, \dots, n$$

or in matrix notation:

$$\mathbf{x}^{(k+1)} = \underline{\mathbf{D}}^{-1} (\mathbf{b} - \underline{\mathbf{L}}\mathbf{x}^{(k)} - \underline{\mathbf{U}}\mathbf{x}^{(k)}),$$

where

$$\underline{\mathbf{A}} = \underline{\mathbf{L}} + \underline{\mathbf{D}} + \underline{\mathbf{U}} = \begin{array}{|c|} \hline \text{0} \\ \hline \end{array} + \begin{array}{|c|} \hline a_{11} & \text{0} \\ \vdots & \vdots \\ \text{0} & a_{nn} \\ \hline \end{array} + \begin{array}{|c|} \hline \text{0} \\ \hline \end{array}$$

Can be improved if always the most recent values of  $\mathbf{x}$  are used:

### 2.2 Gauss–Seidel Method

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$$

or

$$\mathbf{x}^{(k+1)} = \underline{\mathbf{D}}^{-1} (\mathbf{b} - \underline{\mathbf{L}}\mathbf{x}^{(k+1)} - \underline{\mathbf{U}}\mathbf{x}^{(k)}) \iff \mathbf{x}^{(k+1)} = (\underline{\mathbf{D}} + \underline{\mathbf{L}})^{-1} (\mathbf{b} - \underline{\mathbf{U}}\mathbf{x}^{(k)})$$

**Advantages:**

- Averaged faster (often)
- Easier to implement, because  $x_i^{(k)}$  can simply be overwritten by  $x_i^{(k+1)}$

**Diagonal dominance:**

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n$$

Both methods are guaranteed to converge to the correct solution if  $\underline{\mathbf{A}}$  is diagonally dominant.

**Cost:** For each iteration, the cost is  $\mathcal{O}(n^2)$  and only  $\mathcal{O}(n)$  if  $\underline{\mathbf{A}}$  is sparse.

### 3 Polynomial Interpolation

A polynomial  $p : \mathbb{R} \rightarrow \mathbb{R}$  of degree at most  $n$  is most commonly represented in the monomial basis  $\{1, x, x^2, \dots, x^n\}$  as

$$p(x) = \sum_{i=0}^n a_i x^i$$

where  $a_0, \dots, a_n \in \mathbb{R}$ . It can be evaluated in  $\mathcal{O}(n)$  operations using Horner's scheme.

---

**Algorithm 1** Horner's scheme
 

---

**Require:** coefficients  $a_0, \dots, a_n \in \mathbb{R}$ , evaluation point  $x \in \mathbb{R}$

**Ensure:** value  $p(x) = \sum_{i=0}^n a_i x^i$

```

1:  $p \leftarrow a_n$ 
2: for  $i = n - 1, \dots, 0$  do
3:    $p \leftarrow p \cdot x + a_i$ 
4: return  $p$ 

```

---

#### 3.1 Polynomial Curves

In curve design, we consider polynomial curves  $\mathbf{P} : \mathbb{R} \rightarrow \mathbb{R}^d$  (in applications typically  $d \in \{2, 3\}$ , but the theory works for arbitrary  $d$ ). Using the monomial basis, a polynomial curve of degree at most  $n$  can be written as

$$\mathbf{P}(t) = \sum_{i=0}^n \mathbf{A}_i t^i$$

where  $\mathbf{A}_0, \dots, \mathbf{A}_n \in \mathbb{R}^d$ . It can be evaluated with Horner's scheme by replacing scalar operations with vector operations.

---

**Algorithm 2** Horner's scheme
 

---

**Require:** coefficients  $\mathbf{A}_0, \dots, \mathbf{A}_n \in \mathbb{R}^d$ , evaluation point  $t \in \mathbb{R}$

**Ensure:** value  $\mathbf{P}(t) = \sum_{i=0}^n \mathbf{A}_i t^i$

```

1:  $\mathbf{P} \leftarrow \mathbf{A}_n$ 
2: for  $i = n - 1, \dots, 0$  do
3:    $\mathbf{P} \leftarrow \mathbf{P} \cdot t + \mathbf{A}_i$ 
4: return  $\mathbf{P}$ 

```

---

Often we are interested in finding an interpolating polynomial curve. Given points  $\mathbf{P}_0, \dots, \mathbf{P}_n \in \mathbb{R}^d$  and distinct parameter values  $t_0, \dots, t_n \in \mathbb{R}$ , we seek a polynomial curve  $\mathbf{P}$  of degree at most  $n$  such that

$$\mathbf{P}(t_i) = \mathbf{P}_i$$

where  $i = 0, \dots, n$ .

Writing  $\underline{\mathbf{A}} = [\mathbf{A}_0, \dots, \mathbf{A}_n]^\top \in \mathbb{R}^{(n+1) \times d}$  and  $\underline{\mathbf{P}} = [\mathbf{P}_0, \dots, \mathbf{P}_n]^\top \in \mathbb{R}^{(n+1) \times d}$ , these conditions yield the linear system

$$\underline{\mathbf{V}} \underline{\mathbf{A}} = \underline{\mathbf{P}} \tag{3.1}$$

where  $\underline{\mathbf{V}} \in \mathbb{R}^{(n+1) \times (n+1)}$  is the Vandermonde matrix

$$\underline{\mathbf{V}} = [t_i^j]_{i,j=0}^n = \begin{bmatrix} 1 & t_0 & t_0^2 & \cdots & t_0^n \\ 1 & t_1 & t_1^2 & \cdots & t_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_n & t_n^2 & \cdots & t_n^n \end{bmatrix}$$

In principle, one can solve a system like (3.1) in  $\mathcal{O}(n^3)$ , but the special structure of  $\underline{\mathbf{V}}$  can be used to design an  $\mathcal{O}(n^2)$  algorithm. We do not go into the details here, since interpolation can also be done more directly without ever forming the monomial coefficients.

### 3.2 Neville's Algorithm

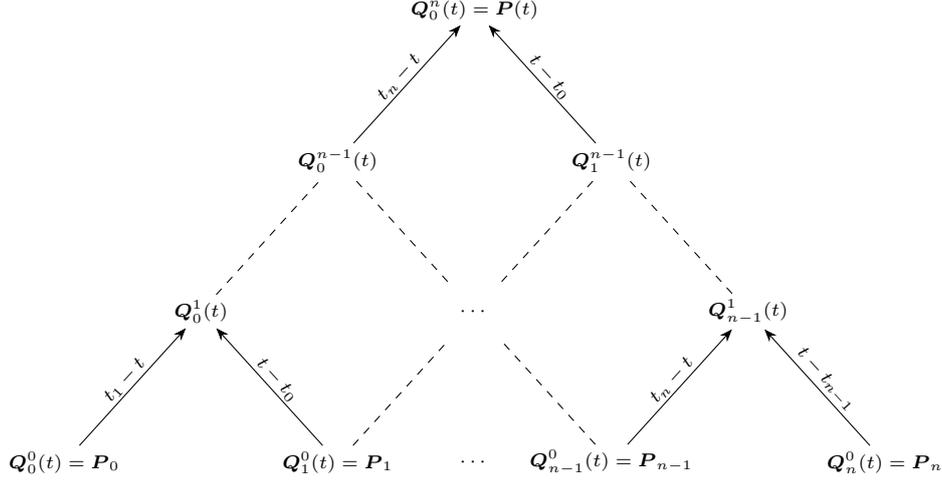
Neville evaluates the interpolating polynomial curve  $P(t)$  at a given  $t \in \mathbb{R}$  *on-the-fly*. It is based on the following recursive formulation. Define the constant curves

$$Q_i^0(t) \equiv P_i$$

for  $i = 0, \dots, n$ . Then define recursively for  $j = 1, \dots, n$ ,

$$Q_i^j(t) = \frac{t_{i+j} - t}{t_{i+j} - t_i} Q_i^{j-1}(t) + \frac{t - t_i}{t_{i+j} - t_i} Q_{i+1}^{j-1}(t)$$

for  $i = 0, \dots, n - j$ . Lower degree curves are linearly blended to form higher degree curves, using affine combinations. The recursion can be visualized as follows:



Note that the common denominators of the affine weights are omitted for clarity. They are always obtained by summing the two weights, i.e.  $t_{i+j} - t_i = (t_{i+j} - t) + (t - t_i)$ .

$Q_0^n(t)$  is the interpolating polynomial curve evaluated at  $t$ . The method costs  $\mathcal{O}(n^2)$  arithmetic operations.

---

#### Algorithm 3 Neville

---

**Require:** data  $(t_i, P_i)$  for  $i = 0, \dots, n$  with distinct  $t_i$ , evaluation parameter  $t \in \mathbb{R}$

**Ensure:** value  $P(t)$  of the interpolating polynomial curve

- 1: **for**  $i = 0, \dots, n$  **do**
  - 2:      $Q_i \leftarrow P_i$  ▷  $Q_i$  stores  $Q_i^j$  in-place
  - 3: **for**  $j = 1, \dots, n$  **do**
  - 4:     **for**  $i = 0, \dots, n - j$  **do**
  - 5:          $Q_i \leftarrow \frac{t_{i+j} - t}{t_{i+j} - t_i} Q_i + \frac{t - t_i}{t_{i+j} - t_i} Q_{i+1}$
  - 6: **return**  $Q_0$
- 

### 3.3 Lagrange Form

From Neville's algorithm, we can derive an explicit representation of the interpolating polynomial, namely the **Lagrange form**

$$P(t) = \sum_{i=0}^n P_i L_i^n(t) \tag{3.2}$$

with Lagrange basis functions

$$L_i^n(t) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}$$

where  $i = 0, \dots, n$ . The interpolation points  $P_0, \dots, P_n$  thus act as coefficients with respect to these Lagrange basis functions of degree  $n$  for parameter values  $t_0, \dots, t_n$ .

Like the monomials, the  $L_i^n$  form a basis of  $\Pi_n$ . They satisfy two important properties:

- Lagrange property:  $L_i^n(t_j) = 0$  for  $j \neq i$  and  $L_i^n(t_i) = 1$
- Partition of unity:  $\sum_{i=0}^n L_i^n(t) \equiv 1$

The first property implies immediately that  $\mathbf{P}(t_i) = \mathbf{P}_i$ .

### 3.4 Barycentric Form

Define  $L(t) = \prod_{i=0}^n (t - t_i)$  and

$$w_i = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{1}{t_i - t_j}$$

where  $i = 0, \dots, n$ . Then the interpolant (3.2) can be written as

$$\mathbf{P}(t) = L(t) \sum_{i=0}^n \frac{w_i}{t - t_i} \mathbf{P}_i$$

and since  $L(t) \sum_{i=0}^n \frac{w_i}{t - t_i} = \sum_{i=0}^n L_i^n(t) \equiv 1$ , as

$$\mathbf{P}(t) = \frac{\sum_{i=0}^n \frac{w_i}{t - t_i} \mathbf{P}_i}{\sum_{i=0}^n \frac{w_i}{t - t_i}} \quad (3.3)$$

This **barycentric form** is typically very stable in floating point arithmetic and yields an  $\mathcal{O}(n)$  evaluation (Algorithm 4) once the weights  $w_i$  are known. To compute them, it takes  $\mathcal{O}(n^2)$  operations. Since they depend only on the nodes  $t_i$ , they can be precomputed and reused for different data  $\mathbf{P}_i$ .

---

**Algorithm 4** Barycentric evaluation of the interpolant

---

**Require:** nodes  $t_0, \dots, t_n$  (distinct), data  $\mathbf{P}_0, \dots, \mathbf{P}_n \in \mathbb{R}^d$ , weights  $w_0, \dots, w_n, t \in \mathbb{R}$

**Ensure:** value  $\mathbf{P}(t)$

```

1:  $\mathbf{N} \leftarrow \mathbf{0}_d, D \leftarrow 0$ 
2: for  $i = 0, \dots, n$  do
3:   if  $t = t_i$  then
4:     return  $\mathbf{P}_i$ 
5:   else
6:      $W \leftarrow \frac{w_i}{t - t_i}$ 
7:      $\mathbf{N} \leftarrow \mathbf{N} + W \mathbf{P}_i$ 
8:      $D \leftarrow D + W$ 
9: return  $\mathbf{N}/D$ 

```

---

### 3.5 Newton Form and Divided Differences

Another representation to write the interpolant (3.2) is the **Newton form**

$$\mathbf{P}(t) = \sum_{i=0}^n \mathbf{B}_i N_i(t) \quad (3.4)$$

with Newton basis functions

$$N_i(t) = \prod_{j=0}^{i-1} (t - t_j)$$

where  $i = 0, \dots, n$ . The Newton basis functions also form a basis of  $\Pi_n$ . The coefficients  $\mathbf{B}_i = \mathbf{P}[t_0, \dots, t_i] \in \mathbb{R}^d$  are the  $i^{\text{th}}$  order *divided differences* of the interpolation points  $\mathbf{P}_i$  with respect to the parameter values  $t_i$  and result from the recursive definition

$$\begin{aligned} \mathbf{P}[t_i] &= \mathbf{P}_i & i &= 0, \dots, n \\ \mathbf{P}[t_i, \dots, t_{i+j}] &= \frac{\mathbf{P}[t_{i+1}, \dots, t_{i+j}] - \mathbf{P}[t_i, \dots, t_{i+j-1}]}{t_{i+j} - t_i} & i &= 0, \dots, n-j, \quad j = 1, \dots, n \end{aligned}$$

As for the Barycentric Form, the coefficients  $\mathbf{B}_i$  can be computed in  $\mathcal{O}(n^2)$  operations as a preprocessing step, and the evaluation then requires  $\mathcal{O}(n)$  operations with a modified Horner scheme (Algorithm 5).

---

**Algorithm 5** modified Horner scheme to evaluate Newton form

---

**Require:** coefficients  $\mathbf{B}_0, \dots, \mathbf{B}_n \in \mathbb{R}^d$ , nodes  $t_0, \dots, t_{n-1}$ , evaluation parameter  $t \in \mathbb{R}$ **Ensure:** value  $\mathbf{P}(t) = \sum_{i=0}^n \mathbf{B}_i N_i(t)$ 

- 1:  $\mathbf{P} \leftarrow \mathbf{B}_n$
  - 2: **for**  $i = n - 1, \dots, 0$  **do**
  - 3:      $\mathbf{P} \leftarrow \mathbf{P} \cdot (t - t_i) + \mathbf{B}_i$
  - 4: **return**  $\mathbf{P}$
-

## 4 Bézier Curves and Splines

Polynomial interpolation is often unsuitable for curve design:

- monomial/Newton coefficients do not correspond intuitively to shape,
- the interpolant can be unstable for (nearly) equidistant parameters,
- changing one data point may affect the curve globally.

These issues motivate representing polynomial curves in the *Bernstein basis*.

### 4.1 Bézier Curves

A Bézier curve  $F : [a, b] \rightarrow \mathbb{R}^d$  of degree  $n$  is defined by control points  $C_0, \dots, C_n \in \mathbb{R}^d$ . It can be evaluated at  $t \in [a, b]$  by repeatedly computing convex combinations (de Casteljau):

$$C_i^0(t) = C_i$$

where  $i = 0, \dots, n$  and for  $j = 1, \dots, n$ ,

$$C_i^j(t) = \frac{b-t}{b-a} C_i^{j-1}(t) + \frac{t-a}{b-a} C_{i+1}^{j-1}(t)$$

where  $i = 0, \dots, n-j$  and then  $F(t) = C_0^n(t)$ .

---

**Algorithm 6** de Casteljau

---

**Require:** control points  $C_0, \dots, C_n \in \mathbb{R}^d$ , interval  $[a, b]$ , parameter  $t \in [a, b]$

**Ensure:** value  $F(t)$

- 1:  $\alpha \leftarrow \frac{b-t}{b-a}, \beta \leftarrow 1 - \alpha$
  - 2: **for**  $i = 0, \dots, n$  **do**
  - 3:      $D_i \leftarrow C_i$
  - 4: **for**  $j = 1, \dots, n$  **do**
  - 5:     **for**  $i = 0, \dots, n-j$  **do**
  - 6:          $D_i \leftarrow \alpha D_i + \beta D_{i+1}$
  - 7: **return**  $D_0$
- 

**Bernstein form.** With the affine reparameterization  $s(t) = \frac{t-a}{b-a} \in [0, 1]$ , the same curve can be written explicitly as

$$F(t) = \sum_{i=0}^n C_i B_i^n(s(t)) \tag{4.1}$$

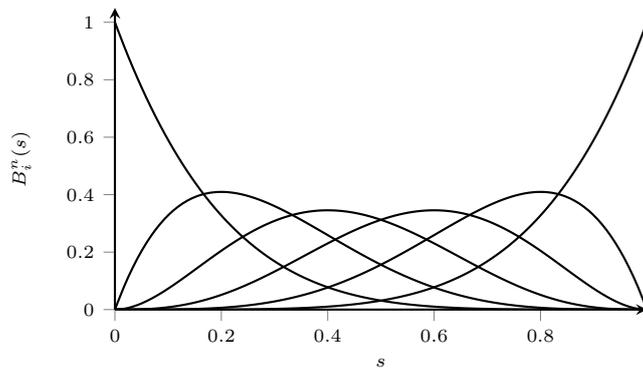
where  $t \in [a, b]$  and the Bernstein polynomials are

$$B_i^n(s) = \binom{n}{i} s^i (1-s)^{n-i} \tag{4.2}$$

where  $i = 0, \dots, n$  and

$$\binom{n}{i} = \frac{n!}{i!(n-i)!}$$

is the binomial coefficient.



**Fact 4.1.** Besides forming a basis for the space of polynomials of degree at most  $n$ , the Bernstein polynomials satisfy the following properties:

- **Partition of unity:**  $\sum_{i=0}^n B_i^n(s) \equiv 1$
- **Positivity:**  $B_i^n(s) \geq 0$  for  $s \in [0, 1]$
- **Endpoints:**  $B_i^n(0) = \delta_{i,0}$  and  $B_i^n(1) = \delta_{i,n}$  (where  $\delta_{i,j}$  is the Kronecker delta)
- **Symmetry:**  $B_i^n(s) = B_{n-i}^n(1-s)$
- **Recursion formula:**  $B_i^n(s) = sB_{i-1}^{n-1}(s) + (1-s)B_i^{n-1}(s)$  with  $B_i^n = 0$  for  $i < 0$  or  $i > n$
- **Derivative:**  $(B_i^n)'(s) = n(B_{i-1}^{n-1}(s) - B_i^{n-1}(s))$
- **Locality:**  $B_i^n$  has a local maximum at  $s = i/n$

These properties almost immediately “translate” to corresponding properties of Bézier curves:

- **Affine invariance:**  $\phi(\sum_{i=0}^n C_i B_i^n(s)) = \sum_{i=0}^n \phi(C_i) B_i^n(s)$  for any affine mapping  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$
- **Convex hull:**  $\mathbf{F}(t) \in \text{conv}\{C_0, \dots, C_n\}$
- **Endpoint interpolation:**  $\mathbf{F}(a) = C_0$  and  $\mathbf{F}(b) = C_n$
- **Derivative:**  $\mathbf{F}'(t) = \frac{n}{b-a} \sum_{i=0}^{n-1} (C_{i+1} - C_i) B_i^{n-1}(s)$
- **Endpoint tangent:**  $\mathbf{F}'(a) = \frac{n}{b-a} (C_1 - C_0)$  and  $\mathbf{F}'(b) = \frac{n}{b-a} (C_n - C_{n-1})$
- **Symmetry:**  $\sum_{i=0}^n C_i B_i^n(s) = \sum_{i=0}^n C_{n-i} B_i^n(1-s)$
- **Local control:**  $C_i$  influences  $\mathbf{F}$  the most at  $t = a + i(b-a)/n$

In  $\mathbb{R}^2$ , Bézier curves satisfy the *variation diminishing property*: any line intersects the curve  $\mathbf{F}$  no more often than it intersects its control polygon. Therefore,  $\mathbf{F}$  cannot exhibit unexpected oscillations. ◀

**Remark 4.1.** Recall the binomial theorem:  $(x+y)^n = \sum_{i=0}^n \binom{n}{i} x^i y^{n-i}$  ◀

Another property is

$$\begin{aligned}
 \sum_{i=0}^n \frac{i}{n} B_i^n(s) &= \sum_{i=0}^n \frac{i}{n} \frac{n!}{i!(n-i)!} s^i (1-s)^{n-i} \\
 &= s \sum_{i=1}^n \frac{(n-1)!}{(i-1)!(n-i)!} s^{i-1} (1-s)^{n-i} \\
 &= s \sum_{i=1}^n \frac{(n-1)!}{(i-1)!((n-1)-(i-1))!} s^{i-1} (1-s)^{(n-1)-(i-1)} \\
 &= s \sum_{j=0}^{n-1} \binom{n-1}{j} s^j (1-s)^{(n-1)-j} \\
 &= s \sum_{j=0}^{n-1} B_j^{n-1}(s) \\
 &= s
 \end{aligned}$$

**Theorem 4.2** (Degree elevation of Bézier curves). Given a Bézier curve

$$\mathbf{F}(t) = \sum_{i=0}^n C_i B_i^n(s) \tag{4.3}$$

of degree  $n$  with control points  $C_0, \dots, C_n \in \mathbb{R}^d$ . Then the same curve admits a representation of degree  $n+1$ ,

$$\mathbf{F}(t) = \sum_{i=0}^{n+1} C_i^* B_i^{n+1}(s) \tag{4.4}$$

with control points  $C_0^* = C_0$ ,  $C_i^* = \frac{i}{n+1} C_{i-1} + \frac{n+1-i}{n+1} C_i$  for  $i = 1, \dots, n$ ,  $C_{n+1}^* = C_n$ . ◀

**Proof.** Using (4.2) and the binomial identities

$$\binom{n}{i} = \frac{n+1-i}{n+1} \binom{n+1}{i} \quad \binom{n}{i} = \frac{i+1}{n+1} \binom{n+1}{i+1}$$

we obtain the two relations

$$(1-s)B_i^n(s) = \frac{n+1-i}{n+1} B_i^{n+1}(s) \quad sB_i^n(s) = \frac{i+1}{n+1} B_{i+1}^{n+1}(s)$$

Hence

$$B_i^n(s) = (1-s)B_i^n(s) + sB_i^n(s) = \frac{n+1-i}{n+1} B_i^{n+1}(s) + \frac{i+1}{n+1} B_{i+1}^{n+1}(s)$$

Insert into (4.3):

$$\mathbf{F}(t) = \sum_{i=0}^n \frac{n+1-i}{n+1} \mathbf{C}_i B_i^{n+1}(s) + \sum_{i=0}^n \frac{i+1}{n+1} \mathbf{C}_i B_{i+1}^{n+1}(s)$$

Reindex the second sum by  $j = i + 1$ :

$$\mathbf{F}(t) = \sum_{i=0}^n \frac{n+1-i}{n+1} \mathbf{C}_i B_i^{n+1}(s) + \sum_{j=1}^{n+1} \frac{j}{n+1} \mathbf{C}_{j-1} B_j^{n+1}(s)$$

Now collect the coefficients of  $B_k^{n+1}(s)$ . For  $k = 0$  and  $k = n + 1$  we obtain  $\mathbf{C}_0^* = \mathbf{C}_0$  and  $\mathbf{C}_{n+1}^* = \mathbf{C}_n$ . For  $k = 1, \dots, n$ , both sums contribute and yield

$$\mathbf{C}_k^* = \frac{k}{n+1} \mathbf{C}_{k-1} + \frac{n+1-k}{n+1} \mathbf{C}_k$$

which proves the claim.  $\square$

## 4.2 Bézier Splines

High-degree Bézier curves are rarely used for designing complex curves for 2 reasons:  $\mathcal{O}(n^2)$  evaluation with Algorithm 6 becomes slower and local control is not very effective for large  $n$ , so the shape becomes difficult to control. Instead, one typically joins several low-degree Bézier segments into a spline.

Let  $\mathbf{F}_0 : [t_0, t_1] \rightarrow \mathbb{R}^d$  and  $\mathbf{F}_1 : [t_1, t_2] \rightarrow \mathbb{R}^d$  be Bézier curves of the same degree  $n$  with control points  $\mathbf{C}_{0,0}, \dots, \mathbf{C}_{0,n}$  and  $\mathbf{C}_{1,0}, \dots, \mathbf{C}_{1,n}$ , respectively.

They are called  **$C^k$ -continuous at  $t_1$**  if  $\mathbf{F}_0(t_1) = \mathbf{F}_1(t_1)$ ,  $\mathbf{F}'_0(t_1) = \mathbf{F}'_1(t_1)$ ,  $\dots$ ,  $\mathbf{F}_0^{(k)}(t_1) = \mathbf{F}_1^{(k)}(t_1)$ .

There is also the weaker notion of *geometric continuity*. The two curves  $\mathbf{F}_0$  and  $\mathbf{F}_1$  are called  **$G^1$ -continuous at  $t_1$**  if  $\mathbf{F}_0(t_1) = \mathbf{F}_1(t_1)$  and  $\mathbf{F}'_0(t_1) \parallel \mathbf{F}'_1(t_1)$ , i.e. the tangent directions agree, but their magnitudes may differ. Clearly,  $C^1$ -continuity implies  $G^1$ -continuity, but not conversely.

**Example 4.2 (Curve design).** In graphic editors, curves are often represented as piecewise cubic Bézier curves. At a join point  $P$ ,

- $C^{-1}$ : the segments do not share the same point  $P$  (gap)
- $C^0$ : the segments share the same point  $P$  (corner allowed)
- $G^1$ : the two Bézier handles at  $P$  are collinear (smooth join)
- $C^1$ : the handles are collinear and of equal length (matched derivatives)

Thus,  $G^1$ -continuity corresponds to a visually smooth join, while  $C^1$ -continuity additionally enforces equal parametrization speed.  $\blacktriangleleft$

### 4.2.1 $C^0$ -continuity

By the endpoint interpolation property from Fact 4.1 we have

$$\begin{aligned} \mathbf{F}_0(t_1) &= \mathbf{C}_{0,n} \\ \mathbf{F}_1(t_1) &= \mathbf{C}_{1,0} \end{aligned}$$

so the 2 curves join continuously at  $t_1$  iff  $\mathbf{C}_{0,n} = \mathbf{C}_{1,0}$ .

### 4.2.2 $C^1$ -continuity

Similarly, by the endpoint tangent property from Fact 4.1, we have

$$\begin{aligned}\mathbf{F}'_0(t_1) &= \frac{n}{t_1 - t_0}(\mathbf{C}_{0,n} - \mathbf{C}_{0,n-1}) \\ \mathbf{F}'_1(t_1) &= \frac{n}{t_2 - t_1}(\mathbf{C}_{1,1} - \mathbf{C}_{1,0})\end{aligned}$$

so the 2 curves join with  $C^1$ -continuity at  $t_1$  iff  $\mathbf{C}_{0,n-1}$ ,  $\mathbf{C}_{0,n} = \mathbf{C}_{1,0}$ ,  $\mathbf{C}_{1,1}$  lie on a common line and are an affine image of the parameter values  $t_0$ ,  $t_1$ ,  $t_2$ , i.e.,

$$\mathbf{C}_{0,n} = \mathbf{C}_{1,0} = (1 - \lambda)\mathbf{C}_{0,n-1} + \lambda\mathbf{C}_{1,1}$$

with  $\lambda = \frac{t_1 - t_0}{t_2 - t_0}$  or, equivalently,  $t_1 = (1 - \lambda)t_0 + \lambda t_2$ .

These 2 observations can be used to construct the simplest kind of Bézier splines, see Example 4.3.

**Example 4.3 (Cubic  $C^1$ -continuous Bézier splines).** Given interpolation points  $\mathbf{P}_0, \dots, \mathbf{P}_m \in \mathbb{R}^d$  at parameters  $t_0 < \dots < t_m$  and tangents  $\mathbf{T}_0, \dots, \mathbf{T}_m \in \mathbb{R}^d$ , define  $m$  cubic Bézier pieces  $\mathbf{F}_i : [t_i, t_{i+1}] \rightarrow \mathbb{R}^d$  by control points

$$\begin{aligned}\mathbf{C}_{i,0} &= \mathbf{P}_i \\ \mathbf{C}_{i,1} &= \mathbf{P}_i + \frac{t_{i+1} - t_i}{3} \mathbf{T}_i \\ \mathbf{C}_{i,2} &= \mathbf{P}_{i+1} - \frac{t_{i+1} - t_i}{3} \mathbf{T}_{i+1} \\ \mathbf{C}_{i,3} &= \mathbf{P}_{i+1}\end{aligned}$$

where  $i = 0, \dots, m - 1$ . Then  $\mathbf{F}$  interpolates both positions and tangents, and is  $C^1$  at the knots. This is called a *cubic Hermite spline*.

**Automatic tangent choices.** If tangents are not given, they can be generated from the interpolation points. For the interior tangents  $i = 1, \dots, m - 1$ , 2 common choices are:

- **FMILL:**  $\mathbf{T}_i = \frac{\mathbf{P}_{i+1} - \mathbf{P}_{i-1}}{t_{i+1} - t_{i-1}}$  (finite difference) ↖ milling devices
- **Bessel:**  $\mathbf{T}_i = \mathbf{Q}'_i(t_i)$ ,  $\mathbf{Q}_i$  is the quadratic interpolant to  $\mathbf{P}_{i-1}$ ,  $\mathbf{P}_i$ ,  $\mathbf{P}_{i+1}$  at  $t_{i-1}$ ,  $t_i$ ,  $t_{i+1}$

Endpoint tangents can be chosen via:

- **Overhauser:**  $\mathbf{T}_0 = \mathbf{Q}'_1(t_0)$ ,  $\mathbf{T}_m = \mathbf{Q}'_{m-1}(t_m)$
- **Natural end:**  $\mathbf{F}''_0(t_0) = \mathbf{F}''_{m-1}(t_m) = \mathbf{0} \Rightarrow \mathbf{C}_{0,1} = \frac{\mathbf{C}_{0,0} + \mathbf{C}_{0,2}}{2}$ ,  $\mathbf{C}_{m-1,2} = \frac{\mathbf{C}_{m-1,1} + \mathbf{C}_{m-1,3}}{2}$

For closed curves ( $\mathbf{P}_0 = \mathbf{P}_m$ ), endpoints are treated like interior points with the two choices above. ◀

### 4.2.3 $C^2$ continuity

The two curves  $\mathbf{F}_0$  and  $\mathbf{F}_1$  join with  $C^2$ -continuity at  $t_1$ , if, in addition to the conditions from Section 4.2.2, there exists a point  $\mathbf{D} \in \mathbb{R}^d$  such that

$$\begin{aligned}\mathbf{C}_{0,n-1} &= (1 - \lambda)\mathbf{C}_{0,n-2} + \lambda\mathbf{D} \\ \mathbf{C}_{1,1} &= (1 - \lambda)\mathbf{D} + \lambda\mathbf{C}_{1,2}\end{aligned}$$

with the same  $\lambda = \frac{t_1 - t_0}{t_2 - t_0}$ . This is called the *A-frame condition* and guarantees that  $\mathbf{C}_{0,n-2}$ ,  $\mathbf{C}_{0,n-1}$ ,  $\mathbf{D}$ , as well as  $\mathbf{D}$ ,  $\mathbf{C}_{1,1}$ ,  $\mathbf{C}_{1,2}$  lie on a common line and are affine images of the parameter values  $t_0$ ,  $t_1$ ,  $t_2$ .

**Example 4.4 (Cubic  $C^2$ -continuous Bézier splines).** Let  $\mathbf{P}_0, \dots, \mathbf{P}_m \in \mathbb{R}^d$  be given at parameters  $t_0 < \dots < t_m$ . Define cubic Bézier pieces  $\mathbf{F}_i : [t_i, t_{i+1}] \rightarrow \mathbb{R}^d$  with  $\mathbf{C}_{i,0} = \mathbf{P}_i$ ,  $\mathbf{C}_{i,3} = \mathbf{P}_{i+1}$ . Choose the two interior control points  $\mathbf{C}_{0,1}$  and  $\mathbf{C}_{0,2}$  (or symmetrically  $\mathbf{C}_{m-1,1}$  and  $\mathbf{C}_{m-1,2}$ ) arbitrarily. Then imposing  $C^1$  continuity and the A-frame condition at all interior knots uniquely determines all remaining control points. Hence a cubic  $C^2$ -continuous Bézier spline interpolating the data has exactly two degrees of freedom. ◀

## 4.3 B-Splines

The difficulty of controlling the shape of a Cubic  $C^2$ -continuous Bézier splines motivates B-splines, which are piecewise polynomial curves with built-in maximal smoothness and intuitive local control.

### 4.3.1 Definition and evaluation

A B-spline curve  $\mathbf{F}$  of degree  $n$  with  $m + 1$  pieces is defined by control points  $\mathbf{D}_0, \dots, \mathbf{D}_{n+m} \in \mathbb{R}^d$  and knots

$$\underbrace{t_0 \leq \dots \leq t_{n-1}}_{n \text{ knots}} \leq \underbrace{t_n \leq \dots \leq t_{n+m+1}}_{m+2 \text{ knots}} \leq \underbrace{t_{n+m+2} \leq \dots \leq t_{2n+m+1}}_{n \text{ knots}}$$

i.e. a non-decreasing *knot vector* of length  $2n + m + 2$ . The curve domain is between the  $m + 2$  “internal” knots, i.e.  $\mathbf{F} : [t_n, t_{n+m+1}] \rightarrow \mathbb{R}^d$ . To evaluate  $\mathbf{F}(t)$  for  $t \in [t_n, t_{n+m+1}]$ , first find  $k$  such that  $t_{n+k} \leq t \leq t_{n+k+1}$ . Then repeatedly compute convex combinations:

$$\mathbf{D}_i^0(t) = \mathbf{D}_i$$

where  $i = k, \dots, k + n$ , and for  $j = 1, \dots, n$ ,

$$\mathbf{D}_i^j(t) = \frac{t_{n+i+1} - t}{t_{n+i+1} - t_{i+j}} \mathbf{D}_i^{j-1}(t) + \frac{t - t_{i+j}}{t_{n+i+1} - t_{i+j}} \mathbf{D}_{i+1}^{j-1}(t)$$

where  $i = k, \dots, k + n - j$  and then finally set  $\mathbf{F}(t) = \mathbf{D}_k^n(t)$ . This is the de Boor algorithm; it is the B-spline analogue of Neville and de Casteljau.

Like Neville’s algorithm and de Casteljau’s algorithm, de Boor’s algorithm costs  $\mathcal{O}(n^2)$ , but since here  $n$  is typically small (often  $2 \leq n \leq 5$ ), this is not an issue.

---

#### Algorithm 7 de Boor

---

**Require:** control points  $\mathbf{D}_0, \dots, \mathbf{D}_{n+m} \in \mathbb{R}^d$ , knots  $t_0, \dots, t_{2n+m+1}$ , parameter  $t \in [t_n, t_{n+m+1}]$

**Ensure:** value  $\mathbf{F}(t)$

```

1:  $k \leftarrow 0$ 
2: while  $t_{n+k+1} < t$  do
3:    $k \leftarrow k + 1$ 
4: for  $i = 0, \dots, n$  do
5:    $\mathbf{E}_i \leftarrow \mathbf{D}_{i+k}$ 
6: for  $j = 1, \dots, n$  do
7:   for  $i = 0, \dots, n - j$  do
8:      $\mathbf{E}_i \leftarrow \frac{t_{n+i+k+1} - t}{t_{n+i+k+1} - t_{i+k+j}} \mathbf{E}_i + \frac{t - t_{i+k+j}}{t_{n+i+k+1} - t_{i+k+j}} \mathbf{E}_{i+1}$ 
9: return  $\mathbf{E}_0$ 
    
```

---

The curve  $\mathbf{F}$  is polynomial of degree  $n$  on each interval  $[t_{n+k}, t_{n+k+1}]$  and  $C^{n-1}$  at interior knots if knots are simple (i.e. not repeated).

### 4.3.2 B-spline basis functions

A B-spline curve can be written explicitly as

$$\mathbf{F}(t) = \sum_{i=0}^{n+m} \mathbf{D}_i N_i^n(t) \quad (4.5)$$

where  $t \in [t_n, t_{n+m+1}]$ .  $N_i^n$ ,  $i = 0, \dots, n + m$ , are the *B-spline basis functions* of degree  $n$  with respect to the knots  $t_0, \dots, t_{2n+m+1}$ .

They are defined recursively as follows. Define the piecewise constant functions

$$N_i^0(t) = \chi_{[t_i, t_{i+1})}(t) := \begin{cases} 1 & t \in [t_i, t_{i+1}) \\ 0 & \text{otherwise} \end{cases}$$

for  $i = 0, \dots, 2n + m$ . Then define recursively for  $j = 1, \dots, n$

$$N_i^j(t) = \frac{t - t_i}{t_{i+j} - t_i} N_i^{j-1}(t) + \frac{t_{i+j+1} - t}{t_{i+j+1} - t_{i+1}} N_{i+1}^{j-1}(t)$$

for  $i = 0, \dots, 2n + m - j$  with the convention that terms with zero denominators are treated as zero.

**Fact 4.3.** The functions  $N_i^j$  are piecewise polynomials of degree  $j$  and  $C^{j-1}$ -continuous. They satisfy

- **Partition of unity:**  $\sum_{i=0}^{2n+m-j} N_i^j(t) = 1$  for  $t \in [t_j, t_{2n-j+m+1}]$
- **Positivity:**  $N_i^j(t) > 0$  for  $t \in (t_i, t_{i+j+1})$

- **Local support:**  $N_i^j(t) = 0$  for  $t \notin [t_i, t_{i+j+1}]$
- **Derivative:**  $(N_i^j)'(t) = \frac{j}{t_{i+j}-t_i} N_i^{j-1}(t) - \frac{j}{t_{i+j+1}-t_{i+1}} N_{i+1}^{j-1}(t)$

These translate to B-spline curve properties. For  $t \in [t_{n+k}, t_{n+k+1}]$

- **Convex hull:**  $\mathbf{F}(t) \in \text{conv}\{\mathbf{D}_k, \dots, \mathbf{D}_{k+n}\}$
- **Local control:**  $\mathbf{F}(t)$  depends only on  $\mathbf{D}_k, \dots, \mathbf{D}_{k+n}$
- **Derivative:**  $\mathbf{F}'(t) = \sum_{i=1}^{n+m} \mathbf{D}'_i N_i^{n-1}(t)$  with  $\mathbf{D}'_i = \frac{n}{t_{i+n}-t_i} (\mathbf{D}_i - \mathbf{D}_{i-1})$
- **Affine invariance:**  $\phi(\sum_{i=0}^{n+m} \mathbf{D}_i N_i^n(t)) = \sum_{i=0}^{n+m} \phi(\mathbf{D}_i) N_i^n(t)$  for any affine  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^d$

As for Bézier curves, B-spline curves also satisfy the variation diminishing property in  $\mathbb{R}^2$ . ◁

### 4.3.3 Interpolation

Given points  $\mathbf{P}_0, \dots, \mathbf{P}_{m+1} \in \mathbb{R}^d$ , a standard approach is to interpolate them at the knots  $t_n, \dots, t_{n+m+1}$ , that is

$$\mathbf{F}(t_{n+i}) = \mathbf{P}_i$$

for  $i = 0, \dots, m+1$ . Using (4.5) this yields the linear system

$$\underline{\mathbf{M}} \mathbf{D} = \mathbf{P}$$

with unknown control points  $\mathbf{D} = [\mathbf{D}_0, \dots, \mathbf{D}_{n+m}]^\top$ , interpolation points  $\mathbf{P} = [\mathbf{P}_0, \dots, \mathbf{P}_{m+1}]^\top$ , and

$$\underline{\mathbf{M}} = \begin{bmatrix} N_0^n(t_n) & N_1^n(t_n) & \cdots & N_{n+m}^n(t_n) \\ N_0^n(t_{n+1}) & N_1^n(t_{n+1}) & \cdots & N_{n+m}^n(t_{n+1}) \\ \vdots & \vdots & \ddots & \vdots \\ N_0^n(t_{n+m+1}) & N_1^n(t_{n+m+1}) & \cdots & N_{n+m}^n(t_{n+m+1}) \end{bmatrix} \in \mathbb{R}^{(m+2) \times (n+m+1)}$$

This system is underdetermined with  $n-1$  degrees of freedom. Hence, additional conditions are needed to get a unique solution. Moreover, due to local support,  $\underline{\mathbf{M}}$  is a band matrix with non-zero entries only near the diagonal.

**Example 4.5 (Cubic B-spline interpolation).** For cubic interpolation,  $n = 3$ , the matrix is tridiagonal up to two additional boundary rows. The two free parameters correspond to the two degrees of freedom in Example 4.4. But they can be fixed more conveniently via end conditions! ☺

Common choices are

- **Natural:**  $\mathbf{F}''(t_3) = \mathbf{F}''(t_{m+4}) = \mathbf{0}$
- **Clamped:**  $\mathbf{F}'(t_3) = \mathbf{T}_0$  and  $\mathbf{F}'(t_{m+4}) = \mathbf{T}_{m+1}$
- **Periodic:**  $\mathbf{F}'(t_3) = \mathbf{F}'(t_{m+4})$  and  $\mathbf{F}''(t_3) = \mathbf{F}''(t_{m+4})$
- **Not-a-knot:**  $\lim_{t \nearrow t_4} \mathbf{F}'''(t) = \lim_{t \searrow t_4} \mathbf{F}'''(t)$  and  $\lim_{t \nearrow t_{m+3}} \mathbf{F}'''(t) = \lim_{t \searrow t_{m+3}} \mathbf{F}'''(t)$

These conditions can be expressed in terms of basis functions and control points and added as the first and last row, turning the system into a square invertible system that uniquely determines the control points.

In the functional setting,  $d = 1$ , the natural interpolating cubic B-spline minimizes the bending energy

$$\|f''\|_2^2 = \int_{t_3}^{t_{m+4}} (f''(t))^2 dt$$

among all interpolating  $C^2$  functions. ◀

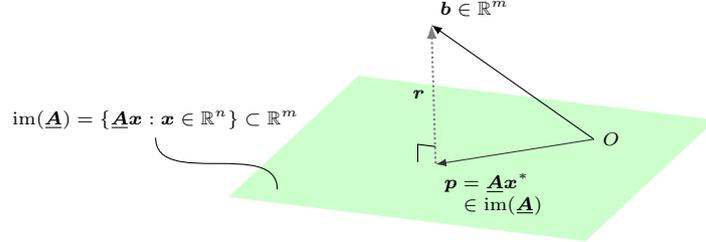
## 5 Linear Least Squares

Let  $\underline{\mathbf{A}} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  and consider the system of equations

$$\underline{\mathbf{A}}\mathbf{x} = \mathbf{b} \quad (5.1)$$

with  $m > n$  (overdetermined system). If  $\mathbf{b}$  does not happen to lie in the hyperplane spanned by the columns of  $\underline{\mathbf{A}}$ , i.e.  $\mathbf{b} \notin \text{im}(\underline{\mathbf{A}})$ , then there is no solution, i.e. it is inconsistent.

For the optimal solution, we require that  $\mathbf{p} = \underline{\mathbf{A}}\mathbf{x}^* \in \text{im}(\underline{\mathbf{A}})$  is as close as possible to  $\mathbf{b}$ , which means that the residual  $\mathbf{r} = \mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*$  is perpendicular to  $\text{im}(\underline{\mathbf{A}})$ :



Note that  $\mathbf{p} \in \text{im}(\underline{\mathbf{A}})$  iff  $\mathbf{p} = \underline{\mathbf{A}}\mathbf{x}^*$  for some  $\mathbf{x}^* \in \mathbb{R}^n$  and that

$$\begin{aligned} \mathbf{b} - \mathbf{p} \text{ perpendicular to } \text{im}(\underline{\mathbf{A}}) &\iff (\mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*) \perp \{\underline{\mathbf{A}}\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\} \\ &\iff (\underline{\mathbf{A}}\mathbf{x})^\top (\mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*) = 0 \quad \forall \mathbf{x} \in \mathbb{R}^n \\ &\iff \mathbf{x}^\top \underline{\mathbf{A}}^\top (\mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*) = 0 \quad \forall \mathbf{x} \in \mathbb{R}^n \\ &\iff \underline{\mathbf{A}}^\top (\mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*) = \mathbf{0}_n \\ &\iff \underline{\mathbf{A}}^\top \mathbf{b} - \underline{\mathbf{A}}^\top \underline{\mathbf{A}}\mathbf{x}^* = \mathbf{0}_n \end{aligned}$$

implying that the least-squares solution satisfies the normal equations

$$\underline{\mathbf{A}}^\top \underline{\mathbf{A}}\mathbf{x}^* = \underline{\mathbf{A}}^\top \mathbf{b} \quad (5.2)$$

The residual  $\mathbf{r} = \mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*$  is usually measured in 3 ways:

- 2-norm:  $\|\mathbf{r}\|_2 = \sqrt{\sum_{i=1}^m r_i^2}$
- squared error:  $\text{SE} = \|\mathbf{r}\|_2^2 = \sum_{i=1}^m r_i^2$
- root mean squared error:  $\text{RMSE} = \sqrt{\frac{1}{m} \sum_{i=1}^m r_i^2} = \frac{\|\mathbf{r}\|_2}{\sqrt{m}}$

### 5.1 Data Fitting (Functional Models)

One important application of least squares is data fitting, where we want to fit a model to given data points.

Suppose we are given  $m$  data points  $(t_1, y_1), \dots, (t_m, y_m)$ , typically coming from measurements or simulations and may contain measurement noise or numerical errors. We want to choose a model  $y = f(t)$  that depends linearly on parameters  $c_1, \dots, c_n$ , i.e.,

$$f(t) = \sum_{j=1}^n c_j \varphi_j(t)$$

for some given basis functions  $\varphi_1, \dots, \varphi_n : \mathbb{R} \rightarrow \mathbb{R}$ . The data fitting workflow can be split into four steps:

1. **Model choice:** Pick a model  $y = f(t)$  with parameters  $\mathbf{c} = (c_1, \dots, c_n)^\top$  that is reasonable from a modeling point of view.
2. **Force the model to match the data.** Plug the data into the model  $y_i = f(t_i) = \sum_{j=1}^n c_j \varphi_j(t_i)$ ,  $i = 1, \dots, m$ , and write this as a linear system

$$\underline{\mathbf{A}}\mathbf{c} = \begin{bmatrix} \varphi_1(t_1) & \dots & \varphi_n(t_1) \\ \vdots & & \vdots \\ \varphi_1(t_m) & \dots & \varphi_n(t_m) \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \mathbf{b}$$

where  $\underline{\mathbf{A}} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$ .

3. **Solve the normal equations:** Find the vector  $\mathbf{c}^*$  that solves  $\underline{\mathbf{A}}^\top \underline{\mathbf{A}} \mathbf{c}^* = \underline{\mathbf{A}}^\top \mathbf{b}$  or use a numerically more stable method (see Section 5.3).
4. **Report an error measure:** Compute the residual  $\mathbf{r} = \mathbf{b} - \underline{\mathbf{A}} \mathbf{c}^*$  and one of the standard error measures.

## 5.2 Curve Fitting (Parametric Models)

The previous setting fits *functions*  $y = f(t)$ . In many geometric applications, however, we want to fit a *parametric curve* in  $\mathbb{R}^d$  to points  $\mathbf{P}_0, \dots, \mathbf{P}_m \in \mathbb{R}^d$ .

1. **Model choice:** Choose basis functions  $B_0, \dots, B_n : \mathbb{R} \rightarrow \mathbb{R}$  (for example monomials, Bernstein polynomials, or B-spline basis functions) and consider the parametric curve

$$\mathbf{F}(t) = \sum_{i=0}^n \mathbf{C}_i B_i(t)$$

where  $\mathbf{C}_i \in \mathbb{R}^d$ . Pick some parameter values  $t_0, \dots, t_m$ .

2. **Force the model to match the data:** Consider the  $m+1$  equations  $\mathbf{P}_k = \mathbf{F}(t_k)$ . Define the matrices

$$\underline{\mathbf{B}} = \begin{bmatrix} B_0(t_0) & \dots & B_n(t_0) \\ \vdots & & \vdots \\ B_0(t_m) & \dots & B_n(t_m) \end{bmatrix} \in \mathbb{R}^{(m+1) \times (n+1)}$$

and  $\underline{\mathbf{C}} = [\mathbf{C}_0, \dots, \mathbf{C}_n]^\top \in \mathbb{R}^{(n+1) \times d}$ ,  $\underline{\mathbf{P}} = [\mathbf{P}_0, \dots, \mathbf{P}_m]^\top \in \mathbb{R}^{(m+1) \times d}$ . Then the fitting conditions  $\mathbf{P}_k \approx \mathbf{F}(t_k)$  can be written compactly as  $\underline{\mathbf{B}} \underline{\mathbf{C}} \approx \underline{\mathbf{P}}$ .

3. **Solve the normal equations:** We now have a matrix-valued least squares problem

$$\min_{\underline{\mathbf{C}}} \|\underline{\mathbf{P}} - \underline{\mathbf{B}} \underline{\mathbf{C}}\|_F^2$$

where  $\|\cdot\|_F$  is the Frobenius norm

$$\|\underline{\mathbf{R}}\|_F^2 = \sum_{i,j} R_{i,j}^2 \quad (5.3)$$

The normal equations become  $\underline{\mathbf{B}}^\top \underline{\mathbf{B}} \underline{\mathbf{C}}^* = \underline{\mathbf{B}}^\top \underline{\mathbf{P}}$  and this decouples into  $d$  standard least squares problems, one for each component of the curve.

4. **Report an error measure:** compute the residual matrix  $\underline{\mathbf{R}} = \underline{\mathbf{P}} - \underline{\mathbf{B}} \underline{\mathbf{C}}^* \in \mathbb{R}^{(m+1) \times d}$  and one of the standard error measures, where the 2-norm is now replaced by the Frobenius norm (5.3).

**Initial parameterization.** We still have to choose parameter values  $t_0, \dots, t_m$ . A common choice is

$$t_0 = 0, \quad t_i = t_{i-1} + \|\mathbf{P}_i - \mathbf{P}_{i-1}\|^\alpha, \quad i = 1, \dots, m,$$

for some exponent  $\alpha \geq 0$ . Typical choices are:

- $\alpha = 0$ : uniform parameterization
- $\alpha = \frac{1}{2}$ : centripetal parameterization
- $\alpha = 1$ : chord length parameterization

In practice, centripetal and chord length parameterizations often yield better fits than uniform parameterization.

**Parameter correction.** Even with a good parameterization, the residual vectors

$$\mathbf{e}_i = \mathbf{P}_i - \mathbf{F}(t_i)$$

usually do not measure the minimal distance between the data point  $\mathbf{P}_i$  and the curve  $\mathbf{F}$ , since the error vectors  $\mathbf{e}_i$  are *not* necessarily orthogonal to the curve. Our goal is now to modify the parameters such that the  $\mathbf{e}_i$  become orthogonal to the curve  $\mathbf{F}$ .

A simple improvement is to locally linearize  $\mathbf{F}$  at  $t_i$  by a first order Taylor expansion

$$\mathbf{F}(t) \approx \mathbf{F}(t_i) + (t - t_i)\mathbf{F}'(t_i) = \mathbf{G}_i(t).$$

We then choose a corrected parameter  $t_i^*$  such that the vector from  $\mathbf{P}_i$  to  $\mathbf{G}_i(t_i^*)$  is orthogonal to the tangent  $\mathbf{F}'(t_i)$ . This leads to the update

$$t_i^* = t_i + \frac{(\mathbf{P}_i - \mathbf{F}(t_i))^\top \mathbf{F}'(t_i)}{\|\mathbf{F}'(t_i)\|^2}$$

Replacing  $t_i$  by  $t_i^*$  (and re-solving the least squares problem) usually reduces the distances  $\|\mathbf{P}_i - \mathbf{F}(t_i)\|$  and improves the overall fit. In practice, one often performs several iterations of curve fitting plus parameter correction.

### 5.3 QR Factorization

The normal equations  $\underline{\mathbf{A}}^\top \underline{\mathbf{A}} \mathbf{x}^* = \underline{\mathbf{A}}^\top \mathbf{b}$  can be numerically problematic, since the condition number of  $\underline{\mathbf{A}}^\top \underline{\mathbf{A}}$  is approximately the square of that of  $\underline{\mathbf{A}}$ . A more sophisticated method is to compute the least-squares solution directly from  $\underline{\mathbf{A}}$  without forming  $\underline{\mathbf{A}}^\top \underline{\mathbf{A}}$ . This can be done via the *QR factorization* of  $\underline{\mathbf{A}}$ .

#### 5.3.1 Classical Gram-Schmidt

Let  $\underline{\mathbf{A}} = (\mathbf{a}_1, \dots, \mathbf{a}_n) \in \mathbb{R}^{m \times n}$  with linearly independent columns. The Gram-Schmidt process constructs orthonormal vectors  $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^m$  spanning the same subspace:

$$\mathbf{y}_j = \mathbf{a}_j - \sum_{k=1}^{j-1} (\mathbf{a}_j^\top \mathbf{q}_k) \mathbf{q}_k, \quad \mathbf{q}_j = \frac{\mathbf{y}_j}{\|\mathbf{y}_j\|}$$

If we set  $r_{k,j} = \mathbf{a}_j^\top \mathbf{q}_k$ ,  $k < j$  and  $r_{j,j} = \|\mathbf{y}_j\|$ , then we can write

$$\underline{\mathbf{A}} = \underline{\mathbf{Q}} \underline{\mathbf{R}}$$

where  $\underline{\mathbf{Q}} = (\mathbf{q}_1, \dots, \mathbf{q}_n) \in \mathbb{R}^{m \times n}$  has orthonormal columns and  $\underline{\mathbf{R}} \in \mathbb{R}^{n \times n}$  is upper triangular. This is the *reduced QR factorization* and can be computed using Algorithm 8.

---

#### Algorithm 8 Classical Gram-Schmidt orthogonalization

---

**Require:** matrix  $\underline{\mathbf{A}} \in \mathbb{R}^{m \times n}$ , where  $\mathbf{A}_j$ ,  $j = 1, \dots, n$  are linearly independent

**Ensure:** reduced QR factorization  $\underline{\mathbf{A}} = \underline{\mathbf{Q}} \underline{\mathbf{R}}$  with  $\underline{\mathbf{Q}} \in \mathbb{R}^{m \times n}$ ,  $\underline{\mathbf{R}} \in \mathbb{R}^{n \times n}$

```

1: for  $j = 1, \dots, n$  do
2:    $\mathbf{y} \leftarrow \mathbf{A}_j$ 
3:   for  $i = 1, \dots, j - 1$  do
4:      $r_{i,j} \leftarrow \mathbf{q}_i^\top \mathbf{A}_j$ 
5:      $\mathbf{y} \leftarrow \mathbf{y} - r_{i,j} \mathbf{q}_i$  ▷ subtract projection of  $\mathbf{A}_j$  onto  $\mathbf{q}_i$ 
6:    $r_{j,j} \leftarrow \|\mathbf{y}\|_2$ 
7:    $\mathbf{q}_j \leftarrow \mathbf{y}/r_{j,j}$  ▷ normalize to unit length

```

---

When the method is successful, it is customary to fill out the matrix of orthogonal unit vectors  $\underline{\mathbf{Q}}$  to a complete basis of  $\mathbb{R}^m$ , to achieve the *full QR factorization*. This can be done, for example, by adding  $m - n$  extra vectors to the  $\mathbf{A}_j$ , so that the  $m$  vectors span  $\mathbb{R}^m$ , and carrying out the Gram-Schmidt method. In terms of the basis of  $\mathbb{R}^m$  formed by  $\mathbf{q}_1, \dots, \mathbf{q}_m$ , the original vectors can be expressed as

$$\underline{\mathbf{A}} = \underline{\mathbf{Q}} \underline{\mathbf{R}}$$

where  $\underline{\mathbf{Q}} \in \mathbb{R}^{m \times m}$  and  $\underline{\mathbf{R}} \in \mathbb{R}^{m \times n}$  with zeros below the  $n \times n$  upper triangular block.

**Definition 5.1.** A square matrix  $\underline{\mathbf{Q}}$  is orthogonal if  $\underline{\mathbf{Q}}^\top = \underline{\mathbf{Q}}^{-1}$  ◀

**Fact 5.1.** An orthogonal matrix  $\underline{Q} \in \mathbb{R}^{m \times m}$  satisfies

- $\underline{Q}^\top$  also orthogonal
- preserves the Euclidian norm of a vector:

$$\|\underline{Q}\mathbf{x}\|_2 = \sqrt{(\underline{Q}\mathbf{x})^\top(\underline{Q}\mathbf{x})} = \sqrt{\mathbf{x}^\top \underline{Q}^\top \underline{Q} \mathbf{x}} = \sqrt{\mathbf{x}^\top \mathbf{x}} = \|\mathbf{x}\|_2 \quad (5.4)$$

for all  $\mathbf{x} \in \mathbb{R}^m$

- product of two orthogonal matrices is orthogonal
- condition number  $\kappa_2(\underline{Q}) = 1$
- $\kappa_2(\underline{A}) = \kappa_2(\underline{Q}\underline{A}) = \kappa_2(\underline{A}\underline{Q})$  ◁

### 5.3.2 Solving least squares problems

For an overdetermined system  $\underline{A}\mathbf{x} = \mathbf{b}$ ,

$$\|\underbrace{\underline{A}\mathbf{x} - \mathbf{b}}_{=: \mathbf{r}}\|_2 = \|\underline{Q}\underline{R}\mathbf{x} - \mathbf{b}\|_2 = \|\underline{Q}(\underline{R}\mathbf{x} - \underline{Q}^\top \mathbf{b})\|_2 \stackrel{(5.4)}{=} \|\underbrace{\underline{R}\mathbf{x} - \underline{Q}^\top \mathbf{b}}_{=: \mathbf{e}}\|_2$$

since  $\underline{Q}$  is orthogonal. Minimizing the norm of  $\mathbf{r} \in \mathbb{R}^m$  is therefore equivalent to minimizing the norm of  $\mathbf{e} \in \mathbb{R}^m$ . Writing out  $\mathbf{e}$  as

$$\mathbf{e} = \begin{bmatrix} e_1 \\ \vdots \\ e_n \\ \hline e_{n+1} \\ \vdots \\ e_m \end{bmatrix} = \begin{bmatrix} r_{1,1} & \cdots & r_{1,n} \\ & \ddots & \vdots \\ & & r_{n,n} \\ \hline 0 & \cdots & 0 \\ \vdots & & \vdots \\ 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} - \begin{bmatrix} d_1 \\ \vdots \\ d_n \\ \hline d_{n+1} \\ \vdots \\ d_m \end{bmatrix} = \underline{R}\mathbf{x} - \mathbf{d}$$

where  $\mathbf{d} = \underline{Q}^\top \mathbf{b}$ , we notice that the last  $m - n$  components of  $\mathbf{e}$  are always  $-d_{n+1}, \dots, -d_m$ , regardless of  $\mathbf{x}$ . So, all we can do to minimize the norm of  $\mathbf{e}$  is to make the first  $n$  components as small as possible. But we can actually make them zero by solving

$$\hat{\underline{R}}\mathbf{x}^* = \hat{\mathbf{d}} \quad (5.5)$$

where  $\hat{\mathbf{d}} = [d_1, \dots, d_n]^\top$  and  $\hat{\underline{R}} \in \mathbb{R}^{n \times n}$  is the upper triangular part of  $\underline{R}$ . Hence,  $\mathbf{x}^*$  of (5.5) is the least-squares solution of  $\underline{A}\mathbf{x} = \mathbf{b}$ , and the squared error of the residual is  $\|\mathbf{r}\|_2^2 = \|\mathbf{e}\|_2^2 = d_{n+1}^2 + \dots + d_m^2$ .

### 5.3.3 Modified Gram–Schmidt and Householder

An improvement in stability can be achieved by replacing  $\mathbf{A}_j$  in line 4 of Algorithm 8 by the current  $\mathbf{y}$  vector, which is mathematically equivalent since

$$r_{i,j} = \mathbf{q}_i^\top \mathbf{y} = \mathbf{q}_i^\top \left( \mathbf{A}_j - \sum_{k=1}^{i-1} r_{k,j} \mathbf{q}_k \right) = \mathbf{q}_i^\top \mathbf{A}_j - \sum_{k=1}^{i-1} r_{k,j} \underbrace{\mathbf{q}_i^\top \mathbf{q}_k}_{=0} = \mathbf{q}_i^\top \mathbf{A}_j$$

For even better stability and efficiency, QR is often computed with Householder reflections.

## 5.4 Weighted Least Squares

While the least squares solution  $\mathbf{x}^*$  minimizes the 2-norm of the residual  $\mathbf{r} = \mathbf{b} - \underline{A}\mathbf{x}$ , that is,

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{r}\|_2 = \arg \min_{\mathbf{x}} \|\mathbf{r}\|_2^2 = \arg \min_{\mathbf{x}} \sum_{i=1}^m |r_i|^2$$

where  $r_i = b_i - (\underline{A}\mathbf{x})_i$ , some of the residual errors  $r_i$  can be rather big, while others are very small. But sometimes we want to emphasize some data points more than others. If we want, for example, to guarantee some error tolerance, then it might be better to reduce the *maximum* error and search for

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{r}\|_\infty = \arg \min_{\mathbf{x}} \max_{i=1, \dots, m} |r_i| \quad (5.6)$$

where  $r_i = b_i - (\underline{\mathbf{A}}\mathbf{x})_i$ . This can be done by solving a *linear programming* problem, but also via *weighted* least squares. The main idea is to give more importance to the data with big residuals when minimizing the norm of the residual and consider a weighted norm instead,

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_{i=1}^m w_i |b_i - (\underline{\mathbf{A}}\mathbf{x})_i|^2 \quad (5.7)$$

given weights  $w_1, \dots, w_m > 0$ . The sum can be written in matrix form as

$$\sum_{i=1}^m w_i |b_i - (\underline{\mathbf{A}}\mathbf{x})_i|^2 = \|\underline{\mathbf{W}}^{1/2}(\mathbf{b} - \underline{\mathbf{A}}\mathbf{x})\|_2^2$$

where  $\underline{\mathbf{W}} = \text{diag}(w_1, \dots, w_m)$ . The corresponding normal equations are

$$\underline{\mathbf{A}}^\top \underline{\mathbf{W}} \underline{\mathbf{A}} \mathbf{x}^* = \underline{\mathbf{A}}^\top \underline{\mathbf{W}} \mathbf{b}$$

If we want to find a solution of (5.6), it is not enough to just solve the weighted least squares problem once (5.7) with some fixed weights  $w_i$ . Instead, we have to iteratively re-weight the data by updating the weights according to

$$w_i^* = w_i |b_i - (\underline{\mathbf{A}}\mathbf{x}^*)_i| \quad (5.8)$$

for  $i = 1, \dots, m$ , thus increasing the importance of the data with big residuals and decreasing the weight for data with small residuals. After a few iterations of solving (5.7) and correcting the weights as in (5.8), we get the optimal solution  $\mathbf{x}^*$  of (5.6). This solution is characterized by the fact that the corresponding residual  $\mathbf{r}^* = \mathbf{b} - \underline{\mathbf{A}}\mathbf{x}^*$  has at least  $n + 1$  alternating extremal components.

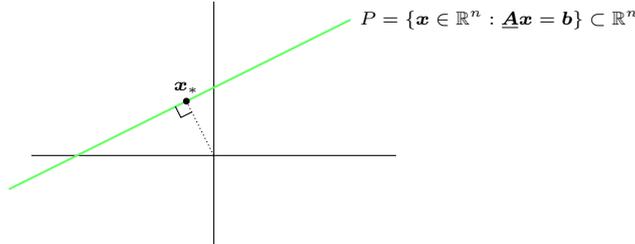
### 5.5 Least-Norm Solutions

If  $\underline{\mathbf{A}} \in \mathbb{R}^{m \times n}$  with  $m < n$ , then the system  $\underline{\mathbf{A}}\mathbf{x} = \mathbf{b}$  is *underdetermined* and typically has infinitely many solutions.

Among all solutions, we are often interested in the one with smallest 2-norm:

$$\mathbf{x}_* = \arg \min \{\|\mathbf{x}\|_2 : \underline{\mathbf{A}}\mathbf{x} = \mathbf{b}\}$$

Since the set  $P = \{\mathbf{x} \in \mathbb{R}^n : \underline{\mathbf{A}}\mathbf{x} = \mathbf{b}\}$  of all solutions can be seen geometrically as an affine subspace (a hyperplane) in  $\mathbb{R}^n$ , this particular least-norm solution  $\mathbf{x}_*$  is the point in  $P$  that is closest to the origin, which is equivalent to saying that  $\mathbf{x}_*$  is orthogonal to  $\mathbf{y} - \mathbf{x}_*$  for any  $\mathbf{y} \in P$ :



Assuming  $\underline{\mathbf{A}}$  has full row rank, then  $\mathbf{x}_*$  is given by

$$\mathbf{x}_* = \underline{\mathbf{A}}^\top (\underline{\mathbf{A}}\underline{\mathbf{A}}^\top)^{-1} \mathbf{b}$$

and it is not hard to verify that this choice satisfies indeed the aforementioned orthogonality condition,

$$\begin{aligned} (\mathbf{y} - \mathbf{x}_*)^\top \mathbf{x}_* &= (\mathbf{y} - \mathbf{x}_*)^\top \underline{\mathbf{A}}^\top (\underline{\mathbf{A}}\underline{\mathbf{A}}^\top)^{-1} \mathbf{b} \\ &= (\underline{\mathbf{A}}(\mathbf{y} - \mathbf{x}_*))^\top (\underline{\mathbf{A}}\underline{\mathbf{A}}^\top)^{-1} \mathbf{b} \\ &= (\underline{\mathbf{A}}\mathbf{y} - \underline{\mathbf{A}}\mathbf{x}_*)^\top (\underline{\mathbf{A}}\underline{\mathbf{A}}^\top)^{-1} \mathbf{b} \\ &= 0 \end{aligned}$$

for all  $\mathbf{y} \in P$ , since  $\mathbf{x}_*, \mathbf{y} \in P$  imply that  $\underline{\mathbf{A}}\mathbf{x}_* = \underline{\mathbf{A}}\mathbf{y} = \mathbf{b}$ . One can view  $\mathbf{x}_*$  as the orthogonal projection of the origin onto the hyperplane  $P$ .

## 5.6 Pseudoinverse

Both the least-squares solution  $\mathbf{x}^*$  and least-norm solution  $\mathbf{x}_*$  can be expressed nicely using the *pseudoinverse*  $\underline{\mathbf{A}}^\dagger \in \mathbb{R}^{n \times m}$  of  $\underline{\mathbf{A}} \in \mathbb{R}^{m \times n}$ :

- If  $m > n$  and  $\underline{\mathbf{A}}$  has full column rank  $n$ , then

$$\underline{\mathbf{A}}^\dagger = (\underline{\mathbf{A}}^\top \underline{\mathbf{A}})^{-1} \underline{\mathbf{A}}^\top$$

and  $\mathbf{x}^* = \underline{\mathbf{A}}^\dagger \mathbf{b}$  is the least-squares solution of  $\underline{\mathbf{A}}\mathbf{x} \approx \mathbf{b}$ .

- If  $m < n$  and  $\underline{\mathbf{A}}$  has full row rank  $m$ , then

$$\underline{\mathbf{A}}^\dagger = \underline{\mathbf{A}}^\top (\underline{\mathbf{A}}\underline{\mathbf{A}}^\top)^{-1}$$

and  $\mathbf{x}_* = \underline{\mathbf{A}}^\dagger \mathbf{b}$  is the least-norm solution of  $\underline{\mathbf{A}}\mathbf{x} = \mathbf{b}$ .

In both cases, if we consider the singular value decomposition

$$\underline{\mathbf{A}} = \underline{\mathbf{U}}\underline{\mathbf{\Sigma}}\underline{\mathbf{V}}^\top$$

of  $\underline{\mathbf{A}}$  and define the pseudoinverse of  $\underline{\mathbf{\Sigma}} \in \mathbb{R}^{m \times n}$  to be the matrix  $\underline{\mathbf{\Sigma}}^\dagger \in \mathbb{R}^{n \times m}$  with the reciprocals  $1/\sigma_i$ ,  $i = 1, \dots, r$ , on the diagonal, where  $r = \min(m, n)$  is the rank of  $\underline{\mathbf{A}}$  and  $\underline{\mathbf{\Sigma}}$ , then

$$\underline{\mathbf{A}}^\dagger = \underline{\mathbf{V}}\underline{\mathbf{\Sigma}}^\dagger\underline{\mathbf{U}}^\top$$

This definition of the pseudoinverse  $\underline{\mathbf{A}}^\dagger$  works even if  $\underline{\mathbf{A}}$  is not a full rank matrix, and in that case  $\mathbf{x} = \underline{\mathbf{A}}^\dagger \mathbf{b}$  is among all vectors that minimize  $\|\mathbf{b} - \underline{\mathbf{A}}\mathbf{x}\|_2$  the one with least norm.

## 6 Nonlinear Least Squares

Suppose we are given  $m$  functions  $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i = 1, \dots, m$  with  $m > n$  and would like to simultaneously solve the  $m$  equations

$$\begin{bmatrix} r_1(\mathbf{x}) \\ \vdots \\ r_m(\mathbf{x}) \end{bmatrix} = \mathbf{0}_m \quad (6.1)$$

in the  $n$  unknowns  $x_1, \dots, x_n$ . As in the case of linear least squares, this is not possible in general, but we can minimize instead the energy  $E : \mathbb{R}^n \rightarrow \mathbb{R}$ , which is given as the sum of squares

$$E(\mathbf{x}) = \sum_{i=1}^m r_i(\mathbf{x})^2 = \mathbf{r}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}) \geq 0 \quad (6.2)$$

where  $\mathbf{r} = [r_1, \dots, r_m]^\top : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . If there exists a point  $\mathbf{x}$  with  $\mathbf{r}(\mathbf{x}) = \mathbf{0}_m$ , then  $E(\mathbf{x}) = 0$  and the minimum is zero; otherwise, any minimizer  $\mathbf{x}^*$  of (6.2) is called a *nonlinear least squares solution* of the system (6.1).

To minimize  $E$ , we set the gradient  $\nabla E : \mathbb{R}^n \rightarrow \mathbb{R}^n$  to zero,

$$\nabla E(\mathbf{x}) = 2 \mathbf{J}_r(\mathbf{x})^\top \mathbf{r}(\mathbf{x}) = \mathbf{0}_n$$

where

$$\mathbf{J}_r(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial r_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial r_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial r_m(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

denotes the Jacobian of  $\mathbf{r}$ . It is convenient to define

$$\mathbf{F}(\mathbf{x}) = \frac{1}{2} \nabla E(\mathbf{x}) = \mathbf{J}_r(\mathbf{x})^\top \mathbf{r}(\mathbf{x}) \in \mathbb{R}^n \quad (6.3)$$

since nonlinear least squares can then be seen as the problem of finding a zero of  $\mathbf{F}$ , which can be done with the multivariate Newton method. Let  $\mathbf{J}_F(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} \mathbf{F}(\mathbf{x}) \in \mathbb{R}^{n \times n}$  denote the Jacobian of  $\mathbf{F}$ . Given a current iterate  $\mathbf{x}^{(k)}$ , a Newton step  $\mathbf{s}^{(k)}$  is obtained by

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left( \mathbf{J}_F(\mathbf{x}^{(k)}) \right)^{-1} \mathbf{F}(\mathbf{x}^{(k)}) \quad (6.4)$$

where, using (6.3) and applying the product rule, we have

$$\begin{aligned} \mathbf{J}_F(\mathbf{x}) &= \frac{\partial}{\partial \mathbf{x}} \left( \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right)^\top \mathbf{r}(\mathbf{x}) \right) \\ &= \left( \frac{\partial}{\partial \mathbf{x}} \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right)^\top \right) \mathbf{r}(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right)^\top \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right) \\ &= \underbrace{\left( \frac{\partial^2}{\partial \mathbf{x}^2} \mathbf{r}(\mathbf{x}) \right)^\top}_{\in \mathbb{R}^{m \times (n \times n)}} \mathbf{r}(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right)^\top \left( \frac{\partial}{\partial \mathbf{x}} \mathbf{r}(\mathbf{x}) \right) \\ &= \underbrace{[\mathbf{H}_{r_1}(\mathbf{x}), \dots, \mathbf{H}_{r_m}(\mathbf{x})]}_{\in \mathbb{R}^{(n \times n) \times m}} \begin{bmatrix} r_1(\mathbf{x}) \\ \vdots \\ r_m(\mathbf{x}) \end{bmatrix} + \mathbf{J}_r(\mathbf{x})^\top \mathbf{J}_r(\mathbf{x}) \\ &= \sum_{i=1}^m r_i(\mathbf{x}) \mathbf{H}_{r_i}(\mathbf{x}) + \mathbf{J}_r(\mathbf{x})^\top \mathbf{J}_r(\mathbf{x}) \end{aligned}$$

where  $\mathbf{H}_{r_i}(\mathbf{x}) \in \mathbb{R}^{n \times n}$  is the Hessian of  $r_i$ . Computing  $\mathbf{J}_F$  exactly involves the calculation of  $mn^2$  second derivatives of the residuals  $r_i$ , which can be expensive. The Gauss-Newton method simplifies Newton's method by replacing  $\mathbf{J}_F$  with the first order approximation

$$\mathbf{J}_F(\mathbf{x}) \approx \mathbf{J}_r(\mathbf{x})^\top \mathbf{J}_r(\mathbf{x}) \quad (6.5)$$

Replacing  $\mathbf{J}_F$  in (6.4) by the approximation (6.5), we obtain a linear least squares problem which can be solved with the same techniques as in the linear case (normal equations, QR factorization, etc.).

**Remark 6.1.** If each residual is linear in  $\mathbf{x}$ , say  $r_i(\mathbf{x}) = b_i - \mathbf{a}_i^\top \mathbf{x}$ , then  $\mathbf{r}(\mathbf{x}) = \mathbf{b} - \mathbf{A}\mathbf{x}$  with  $\mathbf{A}$  having rows  $\mathbf{a}_i^\top$ , and  $\mathbf{J}_r(\mathbf{x}) = -\mathbf{A}$  is constant. In this case, the Gauss-Newton step reduces to the normal equations  $\mathbf{A}^\top \mathbf{A} \mathbf{x}^* = \mathbf{A}^\top \mathbf{b}$ , and the method converges in a single iteration to the linear least-squares solution. ◀

---

**Algorithm 9** Gauss-Newton method

---

**Require:** residual function  $\mathbf{r} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{x}^{(0)}$ ,  $k_{\max}$ ,  $\varepsilon_{\text{tol}} > 0$ **Ensure:** approximation  $\mathbf{x}^*$  that (locally) minimizes  $E(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_2^2$ 

```

1: for  $k = 0, \dots, k_{\max} - 1$  do
2:    $\mathbf{r}^{(k)} \leftarrow \mathbf{r}(\mathbf{x}^{(k)})$ 
3:    $\underline{\mathbf{J}}^{(k)} \leftarrow \underline{\mathbf{J}}_{\mathbf{r}}(\mathbf{x}^{(k)})$ 
4:   solve  $\underline{\mathbf{J}}^{(k)\top} \underline{\mathbf{J}}^{(k)} \mathbf{s}^{(k)} = \underline{\mathbf{J}}^{(k)\top} \mathbf{r}^{(k)}$  for  $\mathbf{s}^{(k)}$ 
5:    $\mathbf{x}^{(k+1)} \leftarrow \mathbf{x}^{(k)} - \mathbf{s}^{(k)}$ 
6:   if  $\|\mathbf{s}^{(k)}\|_2 < \varepsilon_{\text{tol}}$  then
7:     break

```

---

Depending on the initial guess  $\mathbf{x}^{(0)}$ , Algorithm 9 may converge very slowly. This can be improved, either by adding a “regularization term”, which is known as the *Levenberg-Marquardt* method, or by using the exact derivative  $\underline{\mathbf{J}}_{\mathbf{F}}$ .