

## A CONDITION NUMBER THEOREM FOR UNDERDETERMINED POLYNOMIAL SYSTEMS

JÉRÔME DÉGOT

**ABSTRACT.** The condition number of a numerical problem measures the sensitivity of the answer to small changes in the input. In their study of the complexity of Bézout's theorem, M. Shub and S. Smale prove that the condition number of a polynomial system is equal to the inverse of the distance from this polynomial system to the nearest ill-conditioned one. Here we explain how this result can be extended to underdetermined systems of polynomials (that is with less equations than unknowns).

### 1. INTRODUCTION

The condition number of a numerical problem measures the sensitivity of the answer to small changes in the input. We call the problem ill-posed if its condition number is infinite. For a given problem, a condition number theorem asserts that the condition number  $\mu$  is equal to the inverse of the distance of that problem to the set  $\Sigma$  of ill-posed ones.

A first example of such a theorem is due to Eckart and Young [4] about the problem of matrix inversion. See Demmel [3] and Dedieu [2] for a general study concerning condition number theorems for various numerical problems like: matrix inversion, computing eigenvalues and eigenvectors, finding zeroes of polynomials, pole assignment in linear control, . . .

In the case of polynomial systems with the same number of equations as unknowns, Shub and Smale [6] have proved a condition number theorem. We will give a direct and elementary proof of this theorem, using various properties of Bombieri's scalar product. This allows us to extend their result to the case of underdetermined polynomial systems (that is with less equations than unknowns).

### 2. BACKGROUND

Here, we deal only with homogeneous polynomials. Extensions to the affine case are generally trivial, fixing the first variable to 1.

Let  $\mathcal{H}_d$  denote the linear space of all homogeneous polynomial systems  $P = (P_1, \dots, P_m) : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^m$ , where each  $P_i$  is a homogeneous polynomial of  $n + 1$ -variables, of degree  $d_i$  and  $d = (d_1, \dots, d_m)$ .

*Remark 1.* We denote by  $n + 1$  the number of variable, because the first variable  $z_0$  is added to homogenize an ordinary affine polynomial system  $\mathbb{C}^n \rightarrow \mathbb{C}^m$ .

---

Received by the editor August 13, 1996.  
2000 *Mathematics Subject Classification.* Primary 65H10.

For two homogeneous polynomials  $P_i, Q_i : \mathbb{C}^{n+1} \rightarrow \mathbb{C}$  of degree  $d_i$ , Bombieri's scalar product is defined by

$$[P_i, Q_i]_{(d_i)} = \sum_{|\alpha|=d_i} \frac{\alpha!}{d_i!} a_\alpha \overline{b_\alpha},$$

where  $P_i(z) = \sum_{|\alpha|=d_i} a_\alpha z^\alpha$ ,  $Q_i(z) = \sum_{|\alpha|=d_i} b_\alpha z^\alpha$ . With the usual notations, if  $\alpha = (\alpha_0, \dots, \alpha_n) \in \mathbb{N}^{n+1}$ , we have  $|\alpha| = \alpha_0 + \dots + \alpha_n$ ,  $\alpha! = \alpha_0! \dots \alpha_n!$  and  $z^\alpha = z_0^{\alpha_0} \dots z_n^{\alpha_n}$ . This induces a hermitian inner product on  $\mathcal{H}_d$  for  $P, Q \in \mathcal{H}_d$ ,

$$[P, Q] = \sum_{i=1}^m [P_i, Q_i]_{(d_i)},$$

we denote by  $\|P\|$  and  $d(P, Q)$  the associated norm and distance defined by

$$\|P\| = [P, P]^{\frac{1}{2}} \quad \text{and} \quad d(P, Q) = \|P - Q\|.$$

We now give a series of lemmas concerning Bombieri's scalar product. If  $x \in \mathbb{C}^{n+1}$ , we use  $\delta_x$  to denote the homogeneous  $(n+1)$ -variable polynomials of degree 1, defined by

$$\delta_x(z) = \sum_{i=0}^n \overline{x_i} z_i.$$

*Remark 2.* If  $x, y \in \mathbb{C}^{n+1}$ , we have

$$[\delta_x, \delta_y]_{(1)} = \langle y, x \rangle,$$

where  $\langle \cdot, \cdot \rangle$  is the usual inner product of  $\mathbb{C}^{n+1}$ .

The following lemma allows us to replace the evaluation of a polynomial at a point by a duality formula.

**Lemma 1** (B. Reznick [5]). *Let  $P : \mathbb{C}^{n+1} \rightarrow \mathbb{C}$  be a homogeneous polynomial of degree  $d$ . For every  $x \in \mathbb{C}^{n+1}$  we have*

$$P(x) = [P, \delta_x^d]_{(d)}.$$

**Lemma 2** (B. Reznick [5]). *Let  $P, Q$  and  $R$  be three homogeneous polynomials of respective degrees  $p, q$  and  $r$  such that  $p + q = r$ ; then*

$$[PQ, R]_{(r)} = \frac{q!}{r!} [Q, \overline{P}(D)R]_{(q)},$$

where  $P(D)$  is the differential operator defined by

$$P(D) = P\left(\frac{\partial}{\partial x_0}, \dots, \frac{\partial}{\partial x_n}\right).$$

The next lemma is a corollary of Proposition 4 of Beauzamy-Dégot [1].

**Lemma 3.** *Let  $P_1, \dots, P_k$  and  $Q_1, \dots, Q_k$  be homogeneous polynomials of degree 1. We have*

$$[P_1 \dots P_k, Q_1 \dots Q_k]_{(k)} = \frac{1}{k!} \sum_{\sigma} [P_1, Q_{\sigma(1)}] \times \dots \times [P_k, Q_{\sigma(k)}],$$

where  $\sigma$  runs over the set  $S_k$  of all the permutations of  $\{1, \dots, k\}$ .

Combining Lemmas 1 and 2, we obtain a duality formula for polynomial systems.

**Lemma 4.** Let  $P = (P_1, \dots, P_m) : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^m$  be a homogeneous polynomial system, such that each  $P_i$  is of degree  $d_i$ . For  $x, y \in \mathbb{C}^{n+1}$  we have

$$DP(x)y|_i = d_i [P_i, \delta_x^{d_i-1} \delta_y]_{(d_i)} \quad \text{for all } i = 1, \dots, m,$$

where  $DP(x)y|_i$  denotes the  $i$ th coordinate of  $DP(x)y$ , the differential of  $P$  at  $x$  applied to  $y$ .

*Proof.* For all  $i = 1, \dots, m$ , we have

$$DP(x)y|_i = \sum_{j=0}^n y_j \frac{\partial P_i}{\partial x_j}(x).$$

Using Lemmas 1 and 2, we obtain

$$\begin{aligned} DP(x)y|_i &= [\delta_{\bar{y}}(D)P_i, \delta_x^{d_i-1}]_{(d_i-1)} \\ &= d_i [P_i, \delta_x^{d_i-1} \delta_y]_{(d_i)}, \end{aligned}$$

which proves the lemma.

For the proof of the first three lemmas and for a detailed study of the properties of Bombieri's scalar product, we refer the reader to [1].

### 3. THE CONDITION NUMBER OF A POLYNOMIAL SYSTEM

We consider the numerical solution of the equation

$$P(x) = 0$$

where  $P : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^m$  is a homogeneous polynomial system. We want to define the condition number of the system  $P$ . The condition number should measure the relative sensitivity of the solution (output) with respect to the change of the data (input). Let  $\Delta P$  be an infinitesimal perturbation of  $P$ , and let  $\Delta x$  be the first order corresponding perturbation of  $x$ , in the following sense:

$$(1) \quad \|\Delta x\| \simeq \inf\{\|y - x\|; (P + \Delta P)(y) = 0\}.$$

**Lemma 5.** The previous formula implies that  $\Delta x \in (Ker DP(x))^\perp$ .

**Lemma 6.** We have

$$\Delta x = -DP(x)^* \Delta P(x),$$

where  $DP(x)^*$  is the Moore-Penrose inverse of  $DP(x)$ , that is

$$DP(x)^* = (DP(x)|_{(Ker DP(x))^\perp})^{-1}.$$

*Proofs.* The Taylor expansion of  $P + \Delta P$  gives

$$\begin{aligned} 0 &= (P + \Delta P)(x + \Delta x) \\ &= P(x + \Delta x) + \Delta P(x + \Delta x) \\ &= P(x) + DP(x)\Delta x + \Delta p(x) + o(\|\Delta x\|). \end{aligned}$$

Then  $DP(x)\Delta x = -\Delta P(x)$ . By relation (1), we deduce that  $\Delta x \in (Ker DP(x))^\perp$ , and so we have  $\Delta x = -DP(x)^* \Delta P(x)$ , which gives the lemma.

We can now compute the condition number:

$$\begin{aligned} \frac{\|\Delta x\|}{\|x\|} &= \frac{\|DP(x)^* \Delta P(x)\|}{\|x\|} \\ &= \frac{\|DP(x)^* \text{Diag}(\|x\|^{d_i}) \text{Diag}(\|x\|^{-d_i}) \Delta P(x)\|}{\|x\|} \end{aligned}$$

where  $\text{Diag}(\|x\|^{d_i})$  is the diagonal matrix such that

$$\text{Diag}(\|x\|^{d_i}) = \text{Diag}(\|x\|^{d_1}, \dots, \|x\|^{d_m}).$$

Then

$$\frac{\|\Delta x\|}{\|x\|} \leq \frac{\|DP(x)^* \text{Diag}(\|x\|^{d_i})\|_{op}}{\|x\|} \|\text{Diag}(\|x\|^{-d_i}) \Delta P(x)\|_2.$$

By Lemma 1, we have

$$|P_i(x)| = |[P_i, \delta_x^{d_i}]| \leq [P_i] \|x\|^{d_i},$$

so

$$\|\text{Diag}(\|x\|^{-d_i}) \Delta P(x)\|_2 \leq \|\Delta P\|.$$

Therefore,

$$\frac{\|\Delta x\|}{\|x\|} \leq \|DP(x)^* \text{Diag}(\|x\|^{d_i-1})\|_{op} \|\Delta P\|.$$

The condition number of  $P$  at  $x$  is then given by

$$\mu(P, x) = \|DP(x)^* \text{Diag}(\|x\|^{d_i-1})\|_{op}.$$

For sharper estimates on complexity it is convenient to have a further factor of  $\sqrt{d_i}$ . Thus like in Shub-Smale [6], we will use the following normalization of the condition number

$$\mu_{norm}(P, x) = \|DP(x)^* \text{Diag}(\sqrt{d_i} \|x\|^{d_i-1})\|_{op}.$$

#### 4. CONDITION NUMBER THEOREM

It turns out that for many problems of numerical analysis, there is a simple relationship between the condition number of the problem and the shortest distance from that problem to an ill-posed one. A condition number theorem is a theorem which characterizes this relationship. We will give here a general condition number theorem for the problem of finding a zero of a polynomial system of equations.

**Definition 1.** Let  $d = (d_1, \dots, d_m)$  be fixed. We denote by  $\Sigma_x$  the set of all ill-posed problems at  $x$ , that is

$$\Sigma_x = \{Q \in \mathcal{H}_d; Q(x) = 0 \text{ and the rank of } DQ(x) \text{ is less than } m\}.$$

**Theorem 1.** Let  $P : \mathbb{C}^{n+1} \rightarrow \mathbb{C}^m$  be a homogeneous polynomial system with  $n \geq m$ , and let  $x \in \mathbb{C}^{n+1}$  be such that  $P(x) = 0$  and  $\text{Rank}(DP(x)) = m$ . Then

$$\mu_{norm}(P, x) = \frac{1}{d(P, \Sigma_x)}.$$

*Remark 3.* In the case where  $n = m$ , this is the condition number theorem of Shub-Smale [6].

We will give here a proof of this theorem, using the properties of the linear space  $\mathcal{H}_d$  equipped with the hermitian structure described in section 2.

5. PROOF OF THE THEOREM

We have

$$d(P, \Sigma_x) = \inf_{R \in \Sigma_x} \|P - R\| = \inf_{\substack{(P+Q)(x)=0 \\ \text{Rank}(D(P+Q)(x)) < m}} \|Q\|.$$

Then we want to find the polynomial system  $Q \in \mathcal{H}_d$  of smallest norm, such that

(2)  $(P + Q)(x) = 0,$

(3)  $\text{Rank}(D(P + Q)(x)) < m.$

5.1. **First part of the proof:**  $d(P, \Sigma_x) \geq \frac{1}{\mu_{\text{norm}}(P,x)}.$

*Remark 4.* We consider homogeneous polynomial systems, so the Euler identity gives

$$\begin{aligned} (P + Q)(x) &= 0, \\ \iff D(P + Q)(x)x &= 0, \\ \iff x \in \text{Ker}(D(P + Q)(x)). \end{aligned}$$

Let  $Q = (Q_1, \dots, Q_m) \in \mathcal{H}_d$  be one of the homogeneous polynomial systems, with smallest norm satisfying (2) and (3), and let  $x, z_1, \dots, z_k \in \mathbb{C}^{n+1}$ , where  $k \geq n + 1 - m$ , be an orthogonal basis of  $\text{Ker}(D(P + Q)(x))$ . Denote by  $H$  and  $M$  the linear spaces defined by

$$\begin{aligned} M &= \text{Span}\{x, z_1, \dots, z_k\}, \\ H &= \{R \in \mathcal{H}_d; R(x) = 0 \text{ and } DR(x)y = 0 \text{ for all } y \in M\}. \end{aligned}$$

**Proposition 1.** *With the previous notations, we have*

1.  $Q \in H^\perp;$
2.  $H^\perp = \{R \in \mathcal{H}_d; R_i = \delta_x^{d_i-1} \delta_{y_i} \text{ where } y_i \in M\}.$

*Proof.* We can write  $Q = Q_0 + Q_1$ , where  $Q_0 \in H^\perp$  and  $Q_1 \in H$ . Then “ $Q$  satisfies (2) and (3)” implies that “ $Q_0$  satisfies (2) and (3)”. We have

$$\|Q\|^2 = \|Q_0\|^2 + \|Q_1\|^2;$$

thus  $Q_1 = 0$  and then  $Q \in H^\perp$ . For the proof of the second part of the proposition, denote by  $K$  the linear space

$$K = \{R \in \mathcal{H}_d; R_i = \delta_x^{d_i-1} \delta_{y_i} \text{ where } y_i \in M\}.$$

We want to show that  $H^\perp = K$  or equivalently  $H = K^\perp$ . We have

$$Q \in K^\perp \iff [Q_i, \delta_x^{d_i-1} \delta_y]_{(d_i)} = 0 \text{ for all } y \in M \text{ and } i = 1, \dots, m.$$

By Lemma 4, this is equivalent to  $Q(x) = 0$  and  $DQ(x)y = 0$  for all  $y \in M$ , which concludes the proof.

Therefore, the polynomial system  $Q$  is such that

$$Q_i = \delta_x^{d_i-1} \delta_{y_i} \text{ with } y_i \in M,$$

for all  $i = 1, \dots, m$ . Observe first that by Lemmas 1 and 3, we have

$$Q_i(x) = [\delta_x^{d_i-1} \delta_{y_i}, \delta_x^{d_i}]_{(d_i)} = \|x\|^{2(d_i-1)} \langle x, y_i \rangle = 0,$$

then  $\langle x, y_i \rangle = 0$  for all  $i = 1, \dots, m$ . We deduce that

$$\begin{aligned} \|Q\|^2 &= \sum_{i=1}^m [Q_i, Q_i]_{(d_i)} = \sum_{i=1}^m [\delta_x^{d_i-1} \delta_{y_i}, \delta_x^{d_i-1} \delta_{y_i}]_{(d_i)}, \\ (4) \quad \|Q\|^2 &= \sum_{i=1}^m \frac{1}{d_i} \|y_i\|^2 \|x\|^{2(d_i-1)}. \end{aligned}$$

Now we want to estimate  $\|y_i\|$ . By a dimension argument, we know that there exists  $u \in \mathbb{C}^{n+1}$  such that  $\|u\| = 1$  and  $u \in \text{Ker}(DP(x))^\perp \cap \text{Ker}D(P+Q)(x)$ . We then have using Lemma 4,

$$\begin{aligned} D(P+Q)(x)u &= 0, \\ \iff [P_i + Q_i, \delta_x^{d_i-1} \delta_u]_{(d_i)} &= 0 \quad \text{for all } i = 1, \dots, m, \\ \iff [P_i, \delta_x^{d_i-1} \delta_u]_{(d_i)} &= -[\delta_x^{d_i-1} \delta_{y_i}, \delta_x^{d_i-1} \delta_u]_{(d_i)}, \\ &= -\frac{1}{d_i} \|x\|^{2(d_i-1)} \langle u, y_i \rangle, \end{aligned}$$

for all  $i = 1, \dots, m$ .

Then by the Cauchy-Schwarz inequality, we have

$$|[P_i, \delta_x^{d_i-1} \delta_u]_{(d_i)}| \leq \frac{1}{d_i} \|x\|^{2(d_i-1)} \|y_i\| \quad \text{for all } i = 1, \dots, m.$$

Using this last inequality and (4), we have the following chain of calculus:

$$\begin{aligned} \|Q\|^2 &\geq \sum_{i=1}^m \frac{d_i |[P_i, \delta_x^{d_i-1} \delta_u]_{(d_i)}|^2}{\|x\|^{2(d_i-1)}} \\ &= \sum_{i=1}^m \frac{|DP(x)u|_i^2}{d_i \|x\|^{2(d_i-1)}} \\ &= \left\| \text{Diag}\left(\frac{1}{\sqrt{d_i} \|x\|^{d_i-1}}\right) DP(x)u \right\|^2 \\ &\geq \inf_{\substack{u \in \text{Ker}(DP(x))^\perp \\ \|u\|=1}} \left\| \text{Diag}\left(\frac{1}{\sqrt{d_i} \|x\|^{d_i-1}}\right) DP(x)u \right\|^2. \end{aligned}$$

Recall that for an invertible linear operator  $A$ , we have  $\|A^{-1}\|_{op} = \frac{1}{\inf_{\|x\|=1} \|Ax\|}$ , so

$$\begin{aligned} \inf_{\substack{u \in \text{Ker}(DP(x))^\perp \\ \|u\|=1}} \left\| \text{Diag}\left(\frac{1}{\sqrt{d_i} \|x\|^{d_i-1}}\right) DP(x)u \right\| &= \frac{1}{\|DP(x) * \text{Diag}(\sqrt{d_i} \|x\|^{d_i-1})\|_{op}} \\ &= \frac{1}{\mu_{norm}(P, x)}, \end{aligned}$$

which finishes the first part of the proof. To conclude, it remains to prove the opposite inequality.

**5.2. Second part:**  $d(P, x) \leq \frac{1}{\mu_{norm}(P, x)}$ . Assume that  $u \in \text{Ker}(DP(x))^\perp$  is such that  $\|u\| = 1$  and  $u$  realizes the following infimum:

$$\inf_{\substack{y \in \text{Ker}(DP(x))^\perp \\ \|y\|=1}} \left\| \text{Diag}\left(\frac{1}{\sqrt{d_i} \|x\|^{d_i-1}}\right) DP(x)y \right\| = \left\| \text{Diag}\left(\frac{1}{\sqrt{d_i} \|x\|^{d_i-1}}\right) DP(x)u \right\|.$$

Consider the homogeneous polynomial system  $Q = (Q_1, \dots, Q_m) \in \mathcal{H}_d$ , defined by

$$Q_i = -d_i \frac{[P_i, \delta_x^{d_i-1} \delta_u]_{(d_i)}}{\|x\|^{2(d_i-1)}} \delta_x^{d_i-1} \delta_u \quad \text{for all } i = 1, \dots, m.$$

It can easily be seen that

$$\begin{aligned} (P + Q)(x) &= 0, \\ D(P + Q)(x)y &= 0 \quad \text{for all } y \in \text{Ker}(DP(x)), \\ D(P + Q)(x)u &= 0. \end{aligned}$$

Thus  $Q$  satisfies the conditions (2) and (3). Then

$$\begin{aligned} d(P, \Sigma_x) \leq \|Q\| &= \left\| \text{Diag} \left( \frac{1}{\sqrt{d_i} \|x\|^{d_i-1}} \right) DP(x)u \right\| \\ &= \inf_{\substack{y \in \text{Ker}(DP(x))^\perp \\ \|y\|=1}} \left\| \text{Diag} \left( \frac{1}{\sqrt{d_i} \|x\|^{d_i-1}} \right) DP(x)y \right\| \\ &= \frac{1}{\|DP(x)^* \text{Diag}(\sqrt{d_i} \|x\|^{d_i-1})\|_{op}} = \frac{1}{\mu_{norm}(P, x)}, \end{aligned}$$

which ends the proof of the theorem.

## 6. CONCLUSION

It may be surprising that the proof of the theorem does not use what is usually considered as the main property of Bombieri's norm, that is, the unitary invariance (invariance under unitary change of variables in  $P$ ). The use of this property should allow us to replace the point  $x$  at which we compute the condition number of  $P$ , by  $(1, 0, \dots, 0)$ . But this simplification would not improve our proof since we never look at the coefficients of the polynomial system. The second important difference between our approach and that of Shub-Smale, is that we do not use a result on matrices to derive the theorem.

The author would like to thank Mike Shub for helpful discussions.

## REFERENCES

- [1] B. Beauzamy, J. Dégot, *Differential identities*, Trans. Amer. Math. Soc. **347**, (1995), no. 7, pp 2607-2619. MR **96c**:05009
- [2] J-P. Dedieu, *Approximate solutions of numerical problems, condition number analysis and condition number theorems*, The mathematics of numerical analysis. Lectures in Appl. Math., vol. 32, Amer. Math. Soc., Providence, RI, 1996. MR **98a**:65062
- [3] J. Demmel, *On condition numbers and the distance to the nearest ill-posed problem*, Num. Math. **51**, (1987), pp 251-289. MR **88i**:15014
- [4] C. Eckart, G. Young, *The approximation of one matrix by another of lower rank*, Psychometrika **1**, (1936), pp 211-218.
- [5] B. Reznick, *An inequality for products of polynomials*, Proc. Amer. Math. Soc. **117** (1993), 1063-1073. MR **93e**:11058
- [6] M. Shub, S. Smale, *Complexity of Bézout's theorem: geometric aspects*, Journal of the Amer. Math. Soc. **6** (1993), pp 459-501. MR **93k**:65045

LYCÉE FÉNELON, 2, RUE DE L'ÉPERON, 75006 PARIS, FRANCE  
*E-mail address*: jerome.degot@wanadoo.fr