

Tutorial 3

Exercise 1

1) Montrons que $\text{cond}(\sum p_i) = \frac{\sum_{i=1}^n |p_i|}{|\sum_{i=1}^n p_i|}$.

$$\begin{aligned} |\sum_{i=1}^n \tilde{p}_i - \sum_{i=1}^n p_i| &= \left| \sum_{i=1}^n (\tilde{p}_i - p_i) \right| \leq \sum_{i=1}^n |\tilde{p}_i - p_i| \\ &\leq \sum_{i=1}^n 2|p_i| \leq \varepsilon \sum_{i=1}^n |p_i| \end{aligned}$$

on en déduit que :

$$\frac{|\sum_{i=1}^n \tilde{p}_i - \sum_{i=1}^n p_i|}{\varepsilon |\sum_{i=1}^n p_i|} \leq \frac{\varepsilon \sum_{i=1}^n |p_i|}{\varepsilon |\sum_{i=1}^n p_i|} = \frac{\sum_{i=1}^n |p_i|}{|\sum_{i=1}^n p_i|}$$

on a donc $\text{cond}(\sum p_i) \leq \frac{\sum_{i=1}^n |p_i|}{|\sum_{i=1}^n p_i|}$

Si on choisit $\tilde{p}_i = p_i + \varepsilon \text{sign}(p_i) p_i$
 on a $|\tilde{p}_i - p_i| = |\varepsilon \text{sign}(p_i) p_i| = \varepsilon |p_i|$

$$\frac{\left| \sum_{i=1}^n \tilde{p}_i - \sum_{i=1}^n p_i \right|}{\epsilon \left| \sum_{i=1}^n p_i \right|} = \frac{\left| \sum_{i=1}^n \text{sign}(p_i) p_i \right|}{\sum_{i=1}^n \left| \sum_{j=1}^i p_j \right|}$$

$$= \frac{\left| \sum_{i=1}^n \text{sign}(p_i) p_i \right|}{\left| \sum_{i=1}^n p_i \right|} = \frac{\sum_{i=1}^n |p_i|}{\left| \sum_{i=1}^n p_i \right|}$$

En conclusion, $\text{cond}(\sum p_i) = \frac{\sum |p_i|}{\left| \sum p_i \right|}$.

2) Algorithme de sommation récursif

function res = sum(p).

$$\sigma_0 = 0$$

for $i = 1 : n$

$$\sigma_i = \sigma_{i-1} \oplus p_i$$

end

$$res = \sigma_n.$$

Modèle standard : $f(x+y) = x \oplus y$

$$= (1+\delta)(x+y).$$

avec $|S| \leq n$

$$\sigma_0 = 0$$

$$\sigma_1 = p_1$$

$$\sigma_2 = \sigma_1 \oplus p_2 = p_1 \oplus p_2 = (1 + \delta_1)(p_1 + p_2)$$

avec $|\delta_1| \leq \mu$

$$\begin{aligned}\sigma_3 &= \sigma_2 \oplus p_3 = [\sigma_2 + p_3](1 + \delta_2) \text{ avec } |\delta_2| \leq \mu \\ &= [(p_1 + p_2)(1 + \delta_1) + p_3](1 + \delta_2) \\ &= (p_1 + p_2)(1 + \delta_1)(1 + \delta_2) + p_3(1 + \delta_2)\end{aligned}$$

$$\begin{aligned}\sigma_4 &= \sigma_3 \oplus p_4 = [\sigma_3 + p_4](1 + \delta_3) \text{ avec } |\delta_3| \leq \mu \\ &= [(p_1 + p_2)(1 + \delta_1)(1 + \delta_2) + p_3(1 + \delta_2) + p_4] \\ &\quad \times (1 + \delta_3) \\ &= (p_1 + p_2)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) + p_3(1 + \delta_2)(1 + \delta_3) \\ &\quad + p_4(1 + \delta_3)\end{aligned}$$

:

:

:

:

:

:

$$\sigma_n = (\rho_1 + \rho_2) \prod_{i=1}^{n-1} (1 + \delta_i) + \rho_3 \prod_{i=2}^{n-1} (1 + \delta_i) + \rho_4 \prod_{i=3}^{n-1} (1 + \delta_i)$$

$$+ \dots + \rho_n (1 + \delta_{n-1}).$$

On peut montrer que, $|\delta_i| \leq \mu$

$$\prod_{i=1}^n (1 + \delta_i) = 1 + \Theta_n \text{ avec } |\Theta_n| \leq \gamma_n = \frac{n\mu}{1-\mu}$$

$$\sigma_n = (\rho_1 + \rho_2)(1 + \Theta_{n-1}) + \rho_3(1 + \Theta_{n-2}) + \rho_4(1 + \Theta_{n-3})$$

$$+ \dots + \rho_n(1 + \Theta_1).$$

On en déduira que σ_n est la somme exacte
de vecteur $\tilde{\mathbf{p}} = [\underbrace{\rho_1(1 + \Theta_{n-1})}_{\tilde{p}_1}, \underbrace{\rho_2(1 + \Theta_{n-2})}_{\tilde{p}_2}, \underbrace{\rho_3(1 + \Theta_{n-3})}_{\tilde{p}_3}, \dots, \underbrace{\rho_n(1 + \Theta_1)}_{\tilde{p}_n}]$

$$|\tilde{p}_i - p_i| = |\rho_i(1 + \Theta_{n-i+1}) - \rho_i|$$

$$\leq |\Theta_{n-i+1}| |\rho_i| \leq \gamma_{n-i} |\rho_i|$$

L'erreur inverse:

$$\eta(\sigma_n) = \max_{i:1-n} |\frac{\tilde{p}_i - p_i}{p_i}| \leq \gamma_{n-1} \approx (n-1)\mu$$

L'algorithme est bien **inverse-stable**.

3) On cherche une borne pour $|\sigma_n - \sum_{i=1}^n p_i|$

$$\begin{aligned}
 & \left| (\rho_1 + \rho_2)(1 + \theta_{n-1}) + \rho_3(1 + \theta_{n-2}) + \rho_n(1 + \theta_{n-3}) \right. \\
 & \quad \left. + \rho_n(1 + \theta_1) - \sum_{i=1}^n p_i \right| \\
 & \leq |(\rho_1 + \rho_2)\theta_{n-1} + \rho_3\theta_{n-2} + \dots + \rho_n\theta_1| \\
 & \leq |\theta_{n-1}| \sum_{i=1}^n |p_i| \leq \gamma_{n-1} \sum_{i=1}^n |p_i|.
 \end{aligned}$$

Donc $|\sigma_n - \sum_{i=1}^n p_i| \leq \gamma_{n-1} \sum_{i=1}^n |p_i|$

En divisant par $|\sum p_i|$, on obtient

$$\begin{aligned}
 \frac{|\sigma_n - \sum_{i=1}^n p_i|}{|\sum_{i=1}^n p_i|} & \leq \gamma_{n-1} \frac{\sum_{i=1}^n |p_i|}{|\sum_{i=1}^n p_i|} \\
 & \leq \underbrace{\gamma_{n-1}}_{\approx (n-1)\mu} \text{cond}(\sum p_i)
 \end{aligned}$$

4) Mêmes questions avec le produit scalaire.

$$x = (x_i) \in \mathbb{R}^n, y = (y_i) \in \mathbb{R}^n.$$

$$x^T y = \sum_{i=1}^n x_i y_i.$$

* Montrer que $\text{cond}(x^T y) = \frac{2|x|^T|y|}{|x^T y|}$.

Par définition

$$\text{cond}(x^T y) := \limsup_{\varepsilon \rightarrow 0} \left\{ \left| \frac{\tilde{x}^T \tilde{y} - x^T y}{\varepsilon |x^T y|} \right| : \right.$$

$$\begin{aligned} |\tilde{x}_i - x_i| &\leq \varepsilon |x_i| \quad i=1 \dots n \\ |\tilde{y}_i - y_i| &\leq \varepsilon |y_i| \end{aligned}$$

$$\begin{aligned} \left| \frac{\tilde{x}^T \tilde{y} - x^T y}{\varepsilon |x^T y|} \right| &= \left| \frac{(x + \Delta x)^T (y + \Delta y) - x^T y}{\varepsilon |x^T y|} \right| \\ &= \left| \frac{x^T y + (\Delta x)^T y + x^T \Delta y + (\Delta x)^T \Delta y - x^T y}{\varepsilon |x^T y|} \right| \end{aligned}$$

①

$$\leq \frac{|\Delta x|^T |y| + |x|^T |\Delta y| + (\Delta x)^T \Delta y}{\sum |x^T y|}.$$

Or $|\Delta u| \leq \sum |x_i| \text{ and } |\Delta g| \leq \sum |y_j|$

$$\begin{aligned} \text{①} &\leq \frac{\sum |x_i^T y_j| + \sum |x_i^T y_j| + \sum |x_i^T y_j|}{\sum |x^T y|} \\ &\leq \frac{2|x|^T |y|}{|x^T y|} + \sum \frac{|x_i^T y_j|}{|x^T y|}. \end{aligned}$$

$$\Rightarrow \text{cond}(x^T y) \leq \frac{2|x^T y|}{|x^T y|}.$$

Previous

$$\begin{aligned} \tilde{x}_i &= x_i + \sum \text{sign}(y_i) |x_i| \\ \tilde{y}_i &= y_i + \sum \text{sign}(x_i) |y_i| \end{aligned}$$

$$\begin{aligned} \tilde{x}_i \tilde{y}_i &= x_i y_i + \sum x_i \text{sign}(x_i) |y_i| + \sum y_i \text{sign}(y_i) |x_i| \\ &\quad + \sum \text{sign}(y_i) \text{sign}(x_i) |x_i| |y_i| \\ &= x_i y_i + 2 \sum |x_i| |y_i| + \sum \text{sign}(x_i y_i) |x_i| |y_i| \end{aligned}$$

$$\left| \frac{\hat{x}^T \hat{y} - x^T y}{\sum x^T y} \right| = \frac{2 \varepsilon |x^T y| + \varepsilon \sum_{i=1}^n \text{sign}(x_i y_i) |x_i| |y_i|}{|\sum x^T y|}$$

En passant à la limite quand $\varepsilon \rightarrow 0$
 on obtient $\text{cond}(x^T y) = \frac{2 |x^T y|}{|x^T y|}$