

Chapter 6

Pseudo inverse of A :

$$A^+ = (A^*A)^{-1}A^* \quad (6.1)$$

Condition number of $b = Ax$.

$$\kappa(A) = \|A\| \cdot \|A^+\| = \frac{\sigma_{\max}}{\sigma_{\min}} \quad (6.2)$$

Condition number of (A^*A) : (normal equation)

$$\kappa^2(A) \quad (6.3)$$

6.1 Floating Point Arithmetic

$$\beta = \text{radix} \quad (\text{usually } 2) \quad (6.4)$$

$$t = \text{precision} \quad (\text{usually } 24 \text{ or } 53 \text{ in single/double precision}) \quad (6.5)$$

$$x = \pm \left(\frac{m}{\beta^t}\right) \beta^e \quad (6.6)$$

$$m : \text{integer} \quad \beta^{t-1} \leq m \leq \beta^t \quad (6.7)$$

$$e : \text{integer} \quad (6.8)$$

- Machine epsilon:

$$\epsilon_{\text{machine}} = \text{half a unit in the last place} = \frac{1}{2} \frac{1}{\beta^t}. \quad (6.9)$$

- ± 0 (the need for sign of zero)

- Floating point rounding operator

$$\forall x \in \mathbb{R}, \exists \epsilon, |\epsilon| \leq \epsilon_{\text{machine}} : fl(x) = x(1 + \epsilon) \quad (6.10)$$

The distance between x and the closest floating point number is less than $\epsilon_{\text{machine}}$, i.e., less than $\frac{1}{2}$ unit in last place.

- For all practical purposes we say that the result of any floating point operation conforms to:

$$fl(x \odot y) = (x \odot y)(1 + \delta) \quad (6.11)$$

where, $|\delta| \leq \epsilon_{\text{machine}}$

- Infinity ($\pm\infty$)
- Double precision floating point numbers

1	11	52	
sign	exponent	fraction	
	0	0	$\begin{smallmatrix} + \\ - \end{smallmatrix} 0$
	0	$\neq 0$	subnormal
	11...1	0	$\begin{smallmatrix} + \\ - \end{smallmatrix} \text{infinity}$
	11...1	$\neq 0$	NaN

Figure 6.1: FloatingPoint.