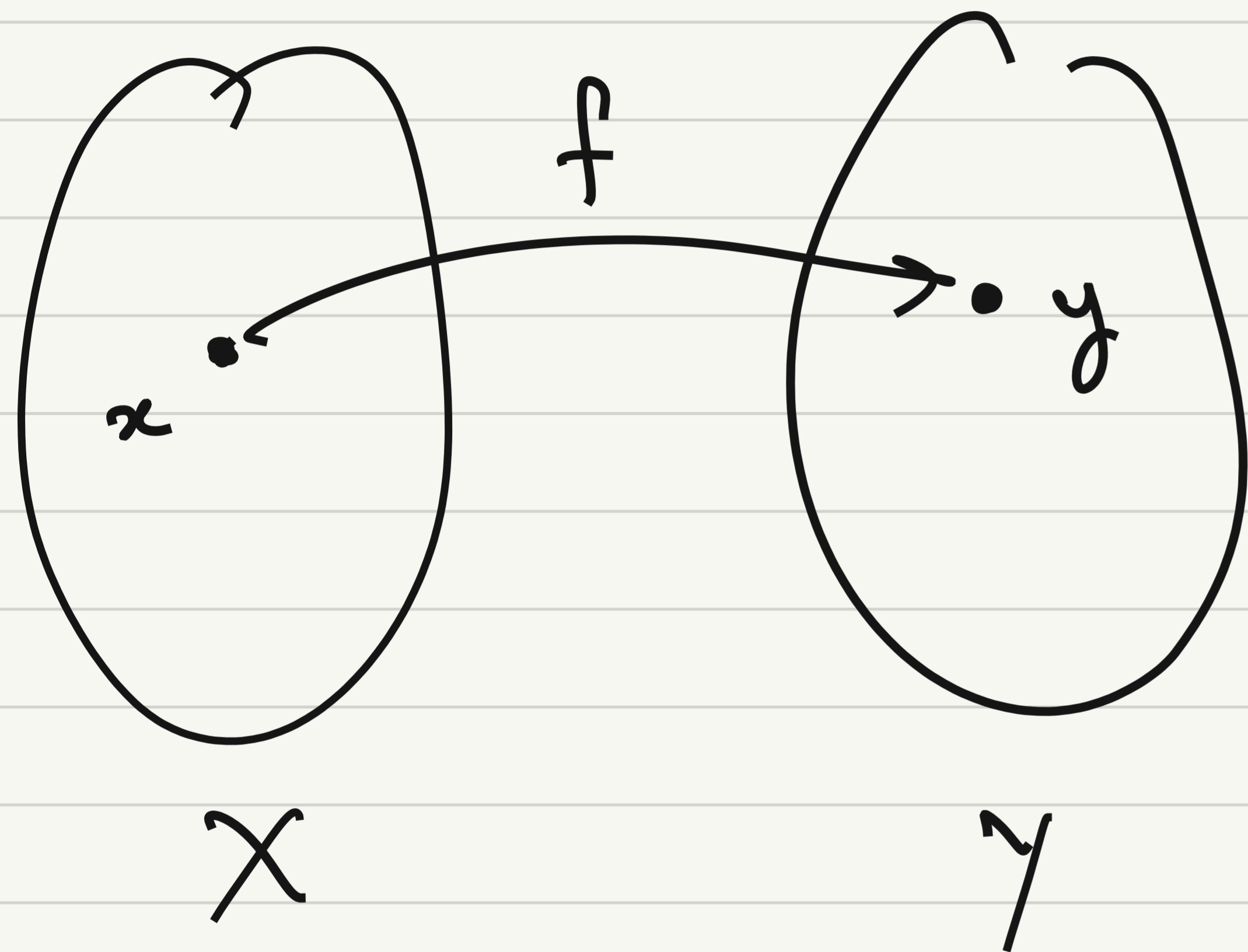


# Floating-Point numbers, Accuracy, Stability

## Review



Condition number:  $\sup_{\delta x} \frac{\text{error in } f(x)}{\text{error in } x}$

$$\hat{\kappa} = \sup_{\delta x} \frac{\|\delta f\|}{\|\delta x\|}$$

$$\kappa = \sup_{\delta x} \frac{\|\delta f\| / \|f\|}{\|\delta x\| / \|x\|} = \frac{\|x\|}{\|f\|} \hat{\kappa}$$

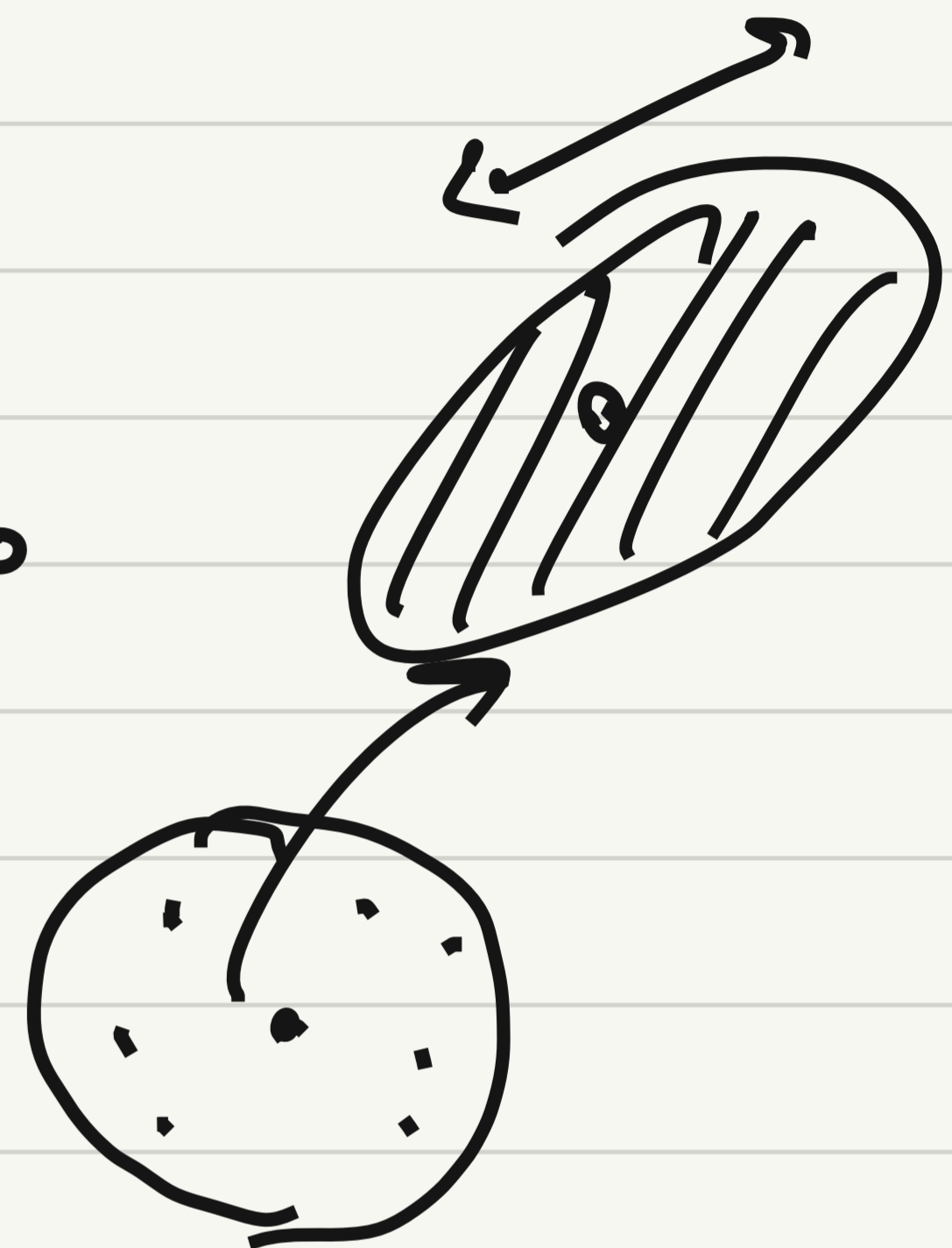
$$f(x_1, x_2) = f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = x_1 - x_2 : \hat{\kappa} = \sup_{\delta x} \frac{2\delta}{\delta} = 2 \quad \text{in } \|\cdot\|_{\infty}$$
$$\kappa = \frac{\|x\|}{\|f\|} \cdot 2 = \frac{\max(|x_1|, |x_2|)}{|x_1 - x_2|} \approx 2$$

roots of polynomials:  $f(a_0, a_1, \dots, a_n) = (z_1, z_2, \dots, z_n)$

$$p(z) = \prod_{i=1}^{20} (z-i) = (z-1)(z-2)\dots(z-20)$$
$$= a_{20}z^{20} + a_{19}z^{19} + \dots + a_1z + a_0$$

$$\kappa \approx 5 \times 10^3$$

$$\frac{\|\delta x\|}{\|x\|} = 10^{-10} \rightarrow \frac{\|\delta f\|}{\|f\|} \leq 5 \times 10^3$$



$$A = \begin{bmatrix} 3 & 0 \\ 0 & 0.5 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$f(x) = Ax = \begin{bmatrix} 3x_1 \\ 0.5x_2 \end{bmatrix}$$

Show that  $\kappa(x) = \frac{3\|x\|}{\|Ax\|}$

whether  $\|x\| = \max(|x_1|, |x_2|)$   
or  $\|x\| = \sqrt{x_1^2 + x_2^2}$

Generalize to  $A = \begin{bmatrix} a_{11} & & \\ & a_{22} & \dots \\ & & a_{nn} \end{bmatrix}$

---

Floating-point arithmetic

IEEE 754:

Single precision: 32 bits (C/C++/Java: float)

Double precision: 64 bits (C/C++/Java: double, Python: float)

Floating point  $\approx$  Scientific notation

$$2022 = \underbrace{2.022}_{\text{mantissa}} \times 10^3$$

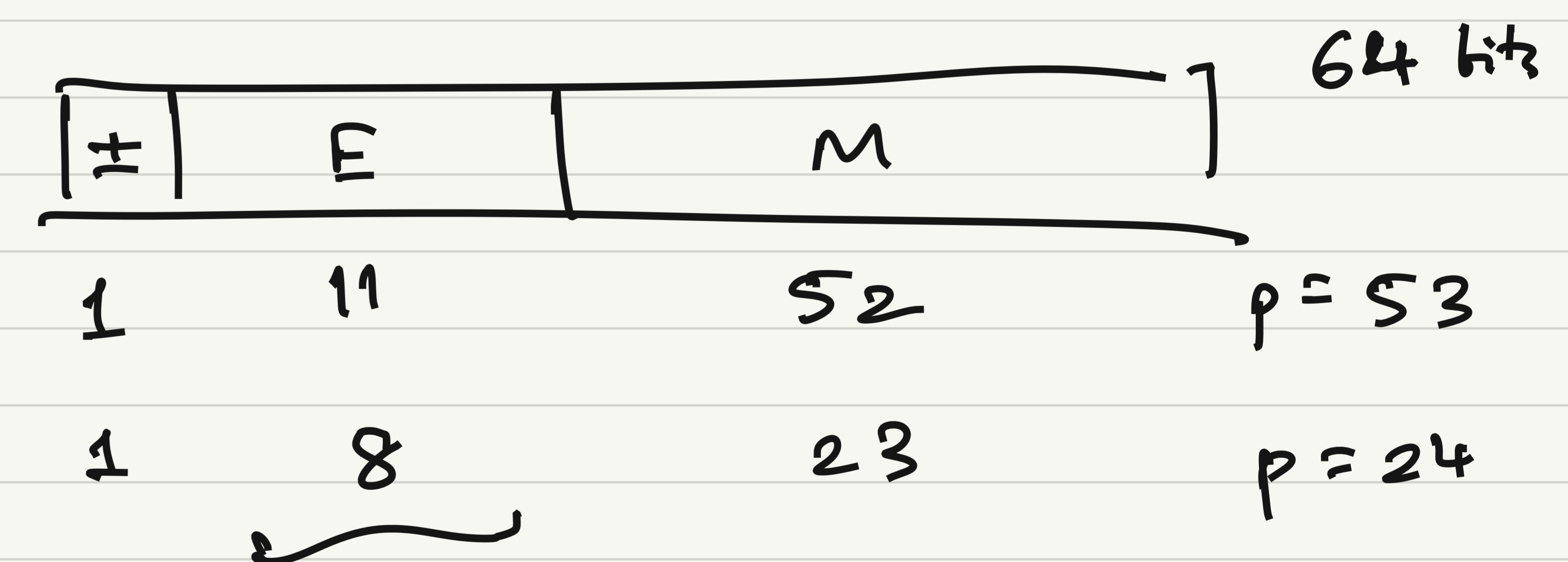
floating point number:  $\pm M \times B^E$  or 0  
 base or **radix** (here 2)

$M = d.dd\dots d$  or  $0.dd\dots d$   
 mantissa  
 p digits

**exponent**

**precision**

double precision  
 ( $\beta = 2$ )



$1.dd\dots d$   
 p digits

$2^{-126} \sim 2^{126}$   
 $10^{-38} \sim 10^{38}$

In our course: Idealized floating point #s

$\beta \in \mathbb{N}$ , fixed  $p$

$M = \underbrace{d \cdot d \cdot d \dots d}_{p \text{ digits}}, E \in \mathbb{Z}$

Goldberg,  
 "What Every CS  
 should know about  
 floats"

$F =$  Set of all floating-point numbers

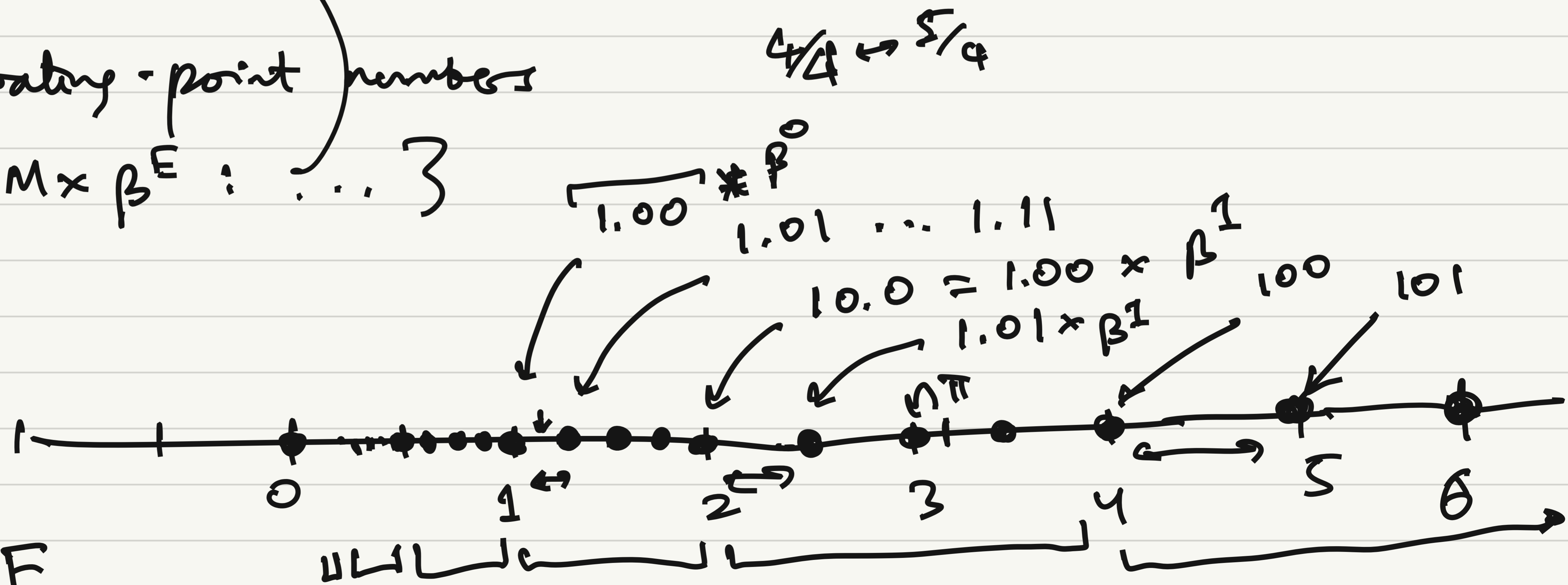
$$= \{0\} \cup \{\pm M \times \beta^E : \dots\}$$

$\beta = 2, p = 3$

$x \in \mathbb{R}$

$fl(x)$

$fl: \mathbb{R} \rightarrow F$



What is the largest relative error between  $x$  and  $fl(x)$ ?

largest rel. error is between 1 and next float =  $1 + \beta^{(1-p)}$

$\beta, p$

$d.ddd\dots d$   
 $\underbrace{\hspace{10em}}_{p \text{ digits}}$

$$x = \frac{1}{2} (1 + (\dots))$$

$$\text{max error} = \frac{1}{2} \beta^{1-p} =: \epsilon_{\text{machine}}$$

1.00...0 ↗  
 1.00...1 ↘

machine epsilon

Single precision

$$\epsilon_m = 2^{-24} \approx 6 \times 10^{-8}$$

Double:  $\epsilon_m = 2^{-53} \approx 10^{-16}$

for any  $x \in \mathbb{R}$ ,  $\underbrace{|fl(x) - x|}_{\leq \epsilon_m |x|} \leq \epsilon_m |x|$

$\Leftrightarrow$   $fl(x) = x(1 + \epsilon)$  for some  $\epsilon$  with  $|\epsilon| \leq \epsilon_m$   
 for all  $x \in \mathbb{R}$

← ①

$$x, y \in F, \quad x+y, x-y, x \cdot y, x/y \notin F$$

$$x \oplus y = fl(x * y)$$

$$\boxed{x \oplus y = (x * y)(1 + \varepsilon) \text{ for some } |\varepsilon| \leq \varepsilon_m} \quad \longleftarrow \textcircled{2}$$

for all  $x, y \in F$

$\varepsilon_m$  : measure of worst-case relative error introduced by floating point system

$x, y, z, \in \mathbb{R}$ , compute  $x - yz$

$$\tilde{x} = fl(x) = x(1 + \epsilon_1) \quad \text{for some } |\epsilon_1| \leq \epsilon_m$$

$$\tilde{y} = fl(y) = y(1 + \epsilon_2)$$

$$\tilde{z} = fl(z) = z(1 + \epsilon_3)$$

$$\tilde{y} \odot \tilde{z} = (\tilde{y} \cdot \tilde{z})(1 + \epsilon_4) = yz(1 + \epsilon_2)(1 + \epsilon_3)(1 + \epsilon_4)$$

$$\tilde{x} \ominus (\tilde{y} \odot \tilde{z}) = (\tilde{x} - \tilde{y} \odot \tilde{z})(1 + \epsilon_5) = \dots$$

$$= \dots = x - yz + \dots$$

$$yz(1 + \underbrace{\epsilon_2 + \epsilon_3 + \epsilon_4}_{\epsilon_6} + \underbrace{\epsilon_2\epsilon_3 + \dots}_{\epsilon_6})$$

$$|\epsilon_2 + \epsilon_3 + \epsilon_4| \leq 3\epsilon_m$$

$$\Rightarrow yz + yz\epsilon_6$$

$$|\epsilon_6| \leq 3\epsilon_m$$



$$\textcircled{1} \quad f(x) = x(1+\varepsilon) \quad \text{for some } |\varepsilon| \leq \varepsilon_m$$

$$\textcircled{2} \quad x \oplus y = (x * y)(1+\varepsilon) \quad \text{for some } |\varepsilon| \leq \varepsilon_m$$

$$0.1 + 0.1 + 0.1 \neq 0.3$$

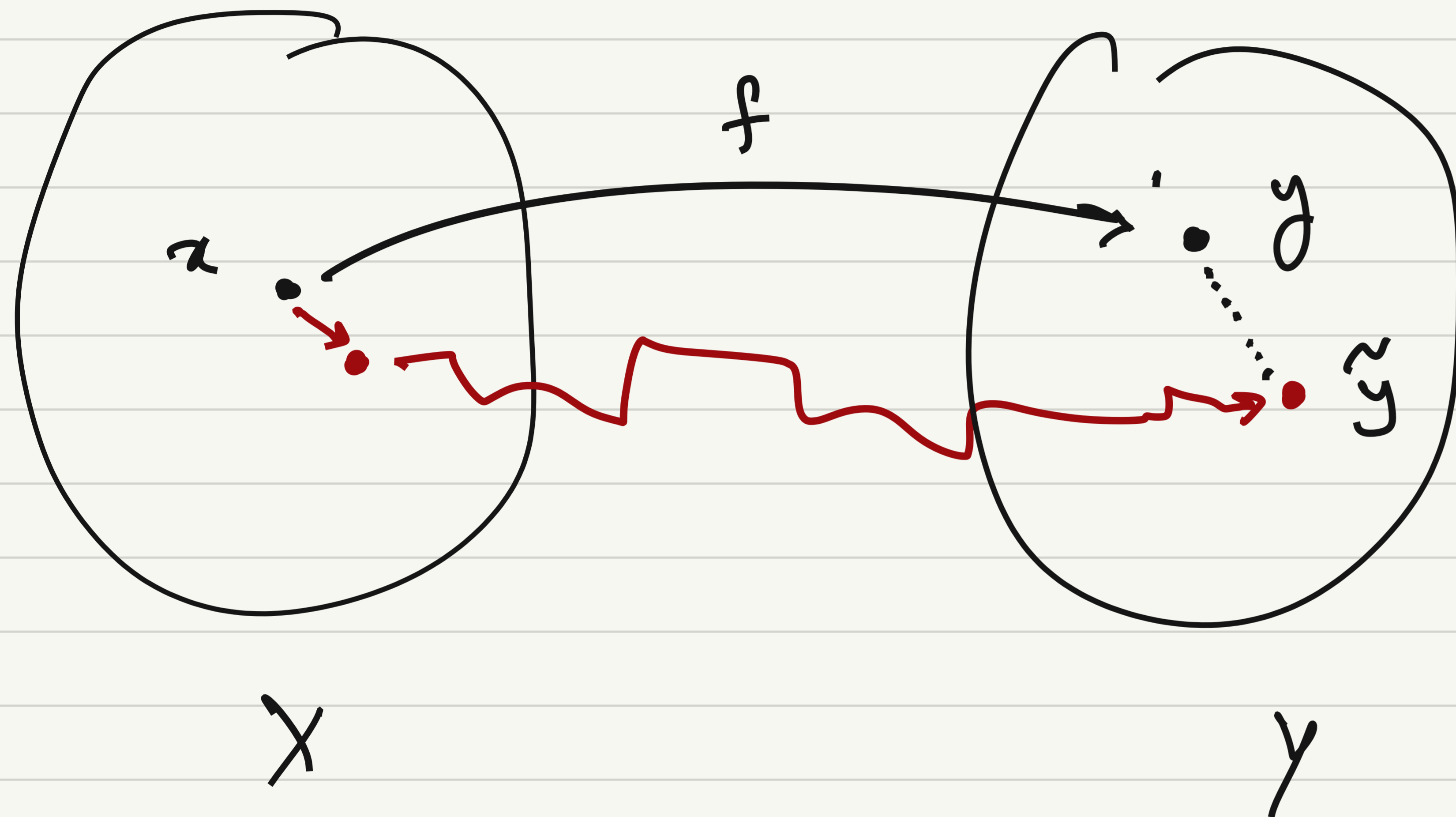
$$a \oplus b = b \oplus a \quad \text{but} \quad a \oplus (b \oplus c) \neq (a \oplus b) \oplus c$$

$$\begin{array}{r} 1.2345\overline{0000} \\ - 1.2321\dots \\ \hline 0.0024 = \underbrace{2.4000}_{\times 10^{-3}} \end{array}$$

$$p(x) = (x-10)^6 \leftarrow$$

$$= x^6 - 6 \times 10 x^5 + 15 \times 10^2 x^4 - \dots + 10^6 \leftarrow$$

# Accuracy and stability of algorithms



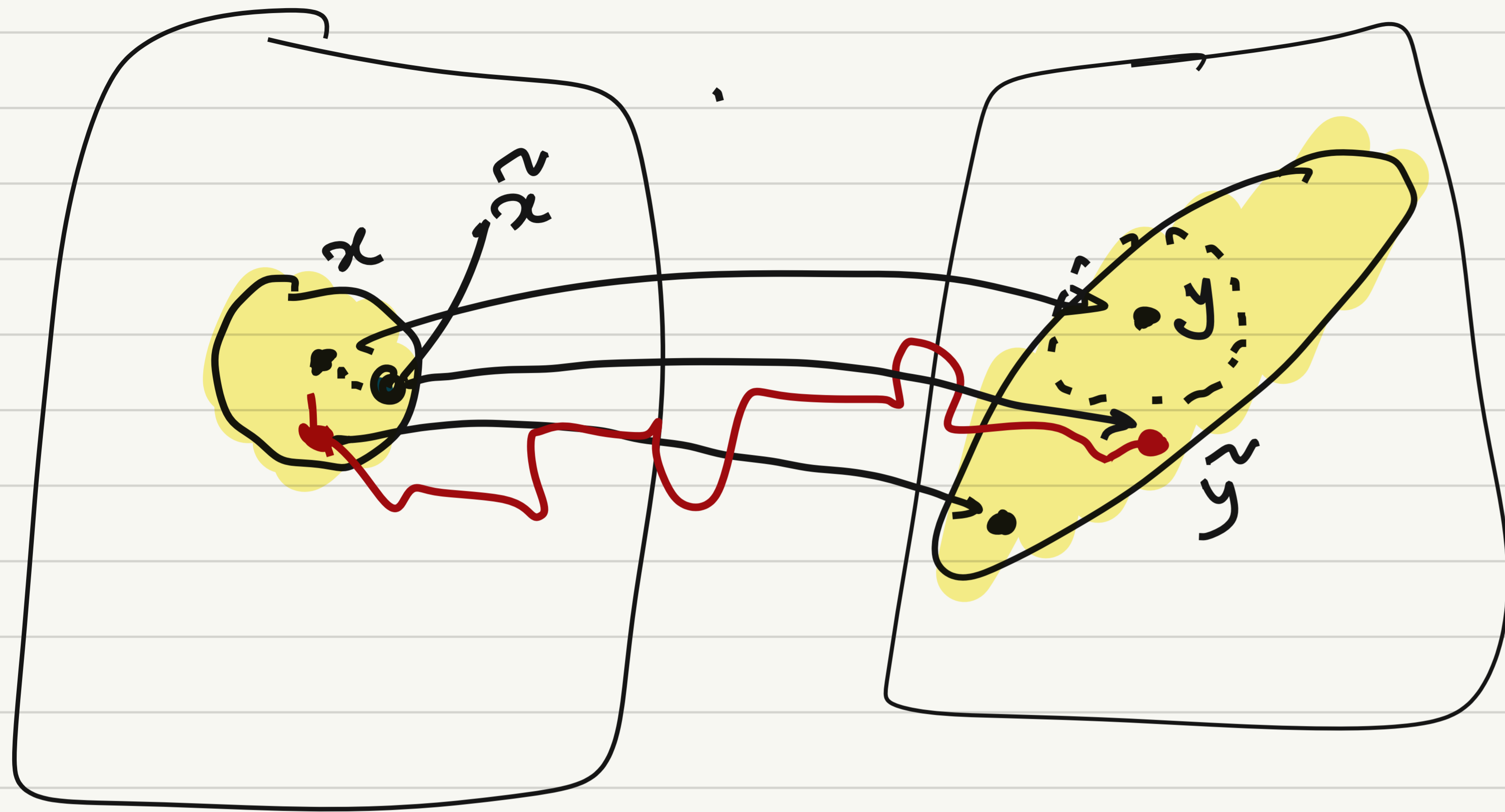
$$\tilde{f}(x) = \tilde{y}$$

$$\text{Abs. forward error} = \|\tilde{y} - y\|$$

$$\text{Rel. " " " " } = \frac{\|\tilde{y} - y\|}{\|y\|}$$

Alg. is accurate if  
rel. f.e. is "small" for all inputs

$$\text{r.f.e.} = \mathcal{O}(\epsilon_m)$$



Computed  $\tilde{y} = \hat{f}(x)$

Suppose there exists  $\tilde{x}$

s.t.  $f(\tilde{x}) = \tilde{y}$

Abs. backward error =  $\|\tilde{x} - x\|$

Rel. b. e. =  $\frac{\|\tilde{x} - x\|}{\|x\|}$

Alg. is backward stable if r.b.e. is small for all  $x$

$$\text{r.b.e.} = O(\epsilon_m)$$

Alg. gives exactly the right answer to nearly the right question

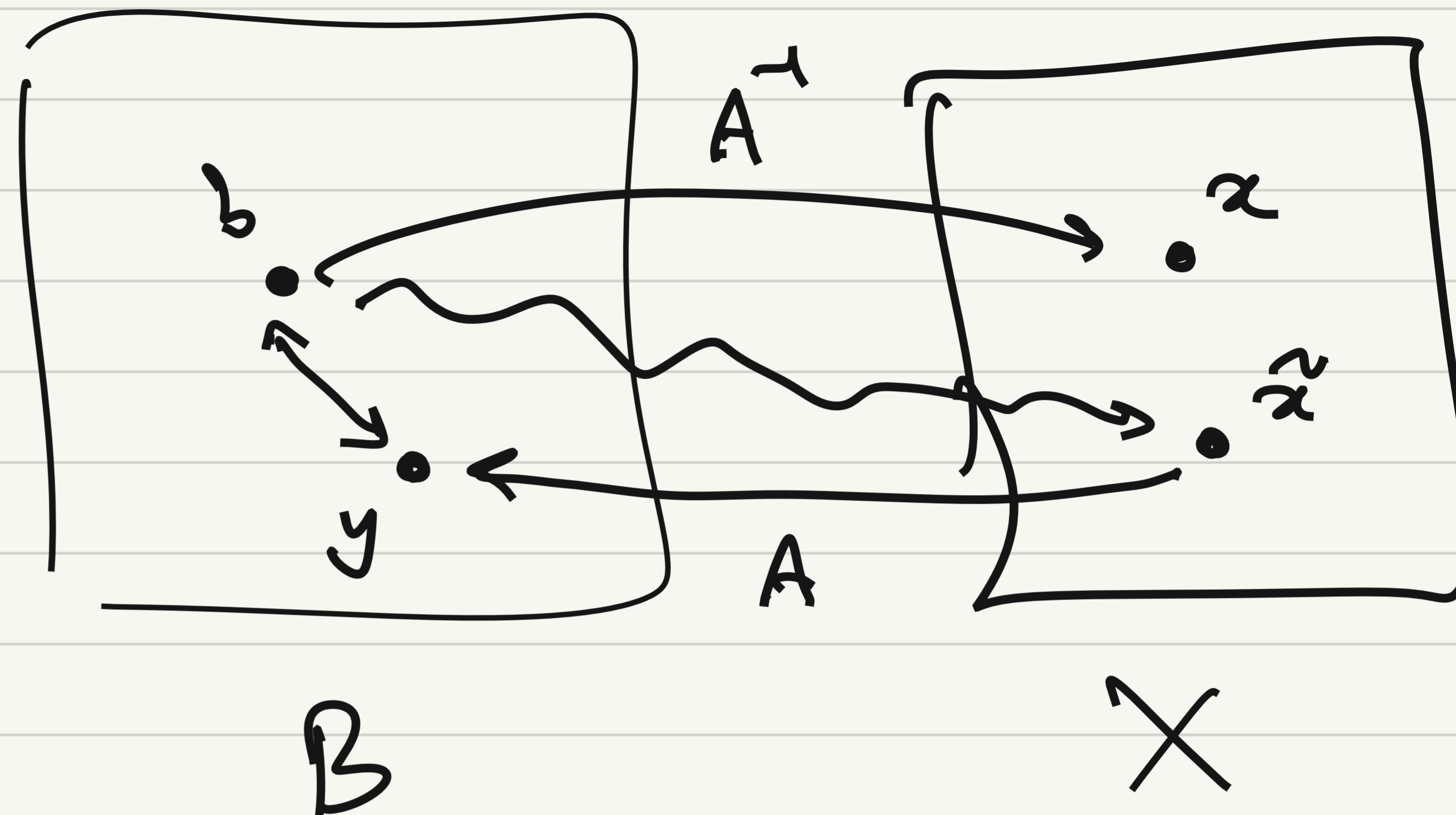
Solve  $Ax = b$ .

$$f(A, b) = A^{-1}b$$

$$\hat{f}(A, b) = \tilde{x}$$

$$A\tilde{x} \rightarrow y$$

$$\|b - A\tilde{x}\|: \text{b.e.}$$



error

$= O(\epsilon_m)$  means  $\exists$  const.  $C$  s.t. for all input  $x$ ,

$$\text{error} \leq C\epsilon_m \quad \text{as } \epsilon_m \rightarrow 0$$