

Calcolo Numerico

Federico Matteoni

A.A. 2019/20

Indice

1	Aritmetica di macchina	5
1.1	Rappresentazione in macchina	5
1.1.1	Aritmetica di Macchina	6
1.1.2	Operazioni di Macchina	7
1.1.3	Errore Inerente	9
1.1.4	Errore Algoritmico	9
1.1.5	Errore Totale	10
1.2	Tecniche per l'Analisi degli Errori	10
1.2.1	Coefficiente di Amplificazione	10
2	Problemi dell'Algebra Lineare Numerica	13
2.1	Norme Matriciali e Vettoriali	13
2.1.1	Norma Vettoriale	13
2.1.2	Norma Matriciale	15
2.1.3	Teoremi di Caratterizzazione	15
2.1.4	Condizionamento della risoluzione di un sistema lineare	16
3	Esercitazioni	17
3.1	Esercitazione 2	17
4	Matlab	19
4.1	Esponenziale	19

Introduzione

Prof.: Luca Gemignani

Calcolo Numerico Metodi numerici per risolvere problemi matematici con il calcolatore. In questo corso ce ne occuperemo dal punto di vista numerico, perché metodi di risoluzione diversi performano in maniera diversa sulla macchina. Cerchiamo di capire quali sono i metodi di interesse e cosa aspettarci dalla loro implementazione.

Il **metodi numerici approssimano la soluzione di problemi matematici**.

Inoltre, il computer **impatta** sul calcolo perché lavora con approssimazioni dei numeri, che su moli elevate di dati e di elaborazioni finiscono per perturbare il risultato ottenuto.

Tipici problemi

$$Ax = b$$

$$Ax = \lambda x$$

$$f(x) = 0$$

$$\int_a^b f(x)dx$$

Matlab Matrix Laboratory, strumento di implementazione per verificare e constatare i risultati teorici.

Informazioni d'esame Compitini, che se complessivamente passati rendono l'orale facoltativo. In alternativa, appelli scritti + orale.

Capitolo 1

Aritmetica di macchina

Modello per capire cosa aspettarci dal punto di vista degli errori dell'esecuzione.

Esempio Per calcolare il limite

$$\lim_{x \rightarrow \infty} \sqrt{x+1} - \sqrt{x}$$

ottengo un caso indeterminato ($\infty - \infty$). Posso semplificare l'espressione ad esempio razionalizzando, e con pochi passaggi ottengo la seguente uguaglianza

$$\sqrt{x+1} - \sqrt{x} = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

Quindi dal punto di vista matematico, le due espressioni sono equivalenti. Ciò **non è sempre vero per la rappresentazione ed il calcolo in macchina**: le due espressioni forniranno risultati completamente diversi. Si rende **necessario**, quindi, **capire quale metodo si comporta meglio** rispetto agli altri.

1.1 Rappresentazione in macchina

Rappresentare i numeri Siamo comunemente abituati a rappresentare un numero in diverse forme.

Ad esempio, $0.1 = \frac{1}{10} = 10 \cdot 0.01$. In generale, per ogni numero ho diversi metodi di rappresentazione \Rightarrow In macchina dobbiamo poter **rappresentare i numeri in maniera univoca**.

Base di numerazione $B \in \mathbb{N}, B > 1$, poiché in base 1 non si può contare.

Una base B ha cifre nell'insieme $\{0, 1, \dots, B-1\}$

Teorema Dato $x \in \mathbb{R}$, con $x \neq 0$

$\exists! (\{d_i\}_{i \geq 1} \text{ con } d_1 \neq 0 \text{ e } d_i \text{ non definitivamente uguali a } B-1) \wedge (p \in \mathbb{Z}) \mid x = \text{segno}(x) \cdot B^p \cdot \sum_{i=1}^{\infty} d_i \cdot B^{-i}$

Considerazioni

Se $x \in \mathbb{C}$ allora viene rappresentato come coppia di numeri reali, quindi il problema si riconduce sempre alla loro rappresentazione

Lo 0 viene rappresentato in modo speciale, poiché non esiste una sua rappresentazione normalizzata

$d_{i \geq 1}$ è una **successione** di cifre

La rappresentazione è **normalizzata** se $d_1 \neq 0$, cioè se la prima cifra è diversa da 0

Non può avere tutte cifre uguali a $B-1$ da un certo punto in poi, la rappresentazione "collassa" al numero successivo

p è detto **esponente**

$\sum_{i=1}^{\infty} d_i \cdot B^{-i}$ è detta **mantissa**

Questa rappresentazione si chiama **in virgola mobile** o **floating point**

Segno	Esponente	Mantissa
-------	-----------	----------

Esempi Ponendo $B = 10$

$$x = 0.01 \Rightarrow x = +10^{-1} \cdot (0.1)$$

$$x = 1.35 \Rightarrow x = +10^1 \cdot (0.135)$$

$$x = 0.0023 \Rightarrow x = +10^{-2} \cdot (0.23)$$

Numeri di Macchina Nei computer ho **registri di lunghezza finita**, quindi essi vengono partizionati: una parte viene riservata alla rappresentazione dell'**esponente** e il resto alla rappresentazione della **mantissa**. L'**insieme dei numeri di macchina** F è quindi così definito

$$F(B, t, m, M) = \left\{ \pm B^p \cdot \sum_{i=1}^t d_i \cdot b^{-i} \mid d_1 \neq 0 \wedge -m \leq p \leq M \right\} \cup \{0\}$$

dove

t sono le **cifre della mantissa**

L'esponente p è compreso tra i valori $-m$ e M

Lo 0 è incluso ma rappresentato a parte

Esempio Ipotizzando di usare registri da 32 bit, posso stanziare 8 bit per l'esponente p (quindi 1 bit per il segno e 7 bit per il valore) e i restanti 24 bit per la mantissa (1 per il segno, 23 per il valore). Quanti numeri posso rappresentare?

$$p \text{ di 7 bit} \Rightarrow 2^7 - 1 = 127 \Rightarrow -127 \leq p \leq 127 \text{ ma lo 0 è rappresentato due volte}$$

$$x = \pm 2^p \sum_{i=1}^2 d_i \cdot 2^{-i}, \text{ con } d_i \in \{0, 1\}, d_1 \neq 0 \Rightarrow d_1 = 1 \text{ sempre}$$

Vedremo che con una serie di accorgimenti è possibile aumentare i numeri esattamente rappresentabili.

Dato quindi $F(B, t, m, M)$, osservo che:

$$F(B, t, m, M) \text{ ha cardinalità finita } N = 2B^{t-1}(B-1)(M+m+1) + 1$$

$$\text{Se } x \in F(B, t, m, M) \wedge x \neq 0 \Rightarrow \omega = B^{-m-1} \leq |x| \leq B^M(1 - B^{-t}) = \Omega$$

Quindi non è possibile rappresentare esattamente numeri non nulli minori di ω . Si può introdurre una rappresentazione **denormalizzata** quando $p = -m$: la condizione $d_1 \neq 0$ può essere abbandonata e posso così rappresentare numeri positivi e negativi compresi in modulo tra B^{-m-t} e $B^{-m}(B^{-1} - B^{-t})$

Analogamente se $p = M$ si introducono rappresentazioni speciali per i simboli $\pm\infty$ e NaN (**not a number**).

1.1.1 Aritmetica di Macchina

La rappresentazione di un numero $x \in R, x \neq 0$ in macchina significa **approssimare** x con un numero $\bar{x} \in F(B, t, m, M)$ commettendo un **errore relativo** di rappresentazione

$$\epsilon_x = \frac{\bar{x} - x}{x} = \frac{\eta_x}{x}, x \neq 0$$

quanto più piccolo possibile in valore assoluto. La quantità $\eta_x = \bar{x} - x$ è detta **errore assoluto** della rappresentazione. L'errore relativo è importante per la **valutazione qualitativa** dell'errore: nelle misurazioni astronomiche un errore assoluto di 1 cm è *qualitativamente diverso* da un errore assoluto di 1 cm nella misurazione di un tavolo.

Dato $x \in R, x \neq 0$, distinguiamo due casi:

1. $|x| < \omega$ (**underflow**) oppure $|x| > \Omega$ (**overflow**)
2. $\omega \leq |x| \leq \Omega$

Nel secondo caso abbiamo quattro tecniche di approssimazione:

1. **Arrotondamento:** x approssimato con il numero rappresentabile \bar{x} più vicino
2. **Troncamento:** x approssimato con il più grande numero rappresentabile \bar{x} tale che $|\bar{x}| \leq |x|$
3. *Round toward $+\infty$:* x approssimato al più piccolo numero rappresentabile maggiore del dato
4. *Round toward $-\infty$:* x approssimato al più grande numero rappresentabile minore del dato

Per semplicità considereremo una macchina che opera per troncamento sull'insieme $F(B, t, m, M)$.

Indicheremo con $trn(x) = \bar{x}$ il risultato dell'approssimazione di x con troncamento e più in generale $fl(x)$ l'**approssimazione in macchina del dato** x nel sistema floating point considerato.

Teorema Sia $x \in R$ con $\omega \leq |x| \leq \Omega$. Si ha

$$|\epsilon_x| = \left| \frac{trn(x) - x}{x} \right| \leq u = B^{1-t}$$

Si osserva che:

$u = B^{1-t}$ è detta **precisione di macchina** ed è **indipendente dalla grandezza del numero**, ma è caratteristica dell'aritmetica floating point, dell'insieme dei numeri rappresentabili e dalla tecnica di approssimazione. Se ad esempio si opera con arrotondamento, u si dimezza.

Per valutare la precisione di macchina possiamo determinare il più piccolo numero di macchina maggiore di 1. Dato x tale numero, abbiamo $x - 1 = |x - 1| = B^{1-t}$ essendo $1 = B^1 \cdot B^{-1}$ rappresentato con esponente $p = 1$. Il seguente script MatLab fornisce il valore richiesto:

```

eps = 0.5;
eps1 = eps + 1;
while(eps > 1)
    eps = 0.5 * eps;
    eps1 = eps + 1;
end
eps = 2 * eps;

```

Dal teorema si ricava che dato $x \in R$, in assenza di overflow/underflow, vale $fl(x) = x(1 + \epsilon_x)$ con $|\epsilon_x| \leq u$.

Questa relazione esprime il modo in cui viene descritto generalmente il legame tra numero reale e sua rappresentazione in macchina.

1.1.2 Operazioni di Macchina

Per le **operazioni aritmetiche** in un sistema floating point si pone un analogo problema di approssimazione, in quanto **il risultato di un'operazione eseguita tra due numeri di macchina in generale non sarà un numero di macchina**.

Operazioni Indichiamo con $\oplus, \ominus, \otimes, \oslash$ le **operazioni aritmetiche di macchina** corrispondenti relativamente all'addizione, sottrazione, prodotto e divisione. Si richiede che le operazioni siano interne all'insieme dei numeri di macchina. Una ragionevole definizione, derivata dal teorema precedente e in assenza di overflow/underflow, risulta:

$$\forall a, b \in F(B, t, m, M), a \oplus b = fl(a + b) = (a + b)(1 + \epsilon_1), |\epsilon_1| \leq u$$

con ϵ_1 detto **errore locale dell'operazione**. Sempre in assenza di overflow/underflow, se $a, b \in R$ si ha

$$fl(a + b) = fl(a) \oplus fl(b) = (a(1 + \epsilon_a) + b(1 + \epsilon_b))(1 + \epsilon_1) \doteq (a + b) + a\epsilon_a + b\epsilon_b + (a + b)\epsilon_1$$

dove con \doteq si indica che l'eguaglianza vale **considerando le sole componenti lineari negli errori** e trascurando le componenti di ordine superiore (**analisi al primo ordine dell'errore**), in virtù del fatto che gli ϵ sono quantità piccole $0 < \epsilon < 1$, quindi trascurabili negli ordini superiori al primo.

Si ottiene che, in assenza di overflow/underflow, se $a, b \in R, a + b \neq 0$, allora

$$\frac{fl(a + b) - (a + b)}{a + b} \doteq \frac{a}{a + b}\epsilon_a + \frac{b}{a + b}\epsilon_b + \epsilon_1$$

che esprime la dipendenza dell'errore totale commesso nel calcolo della somma tra due numeri reali rispetto agli errori generati dall'approssimazione dei dati iniziali (**errore inerente**) e agli errori generati dall'algoritmo di calcolo (**errore algoritmico**) visto come sequenza di operazioni aritmetiche.

Esempio Analizziamo cosa succede in macchina quando proviamo a calcolare $f(x) = x^2 - 1 = (x - 1)(x + 1)$. La **prima situazione di errore** sia ha sulla rappresentazione di x che, **in generale, non è un numero di macchina**.

$$x \rightarrow \tilde{x} = x(1 + \epsilon_x)$$

con $|\epsilon_x| \leq u$. Inoltre, sempre in generale, **le operazioni aritmetiche non sono operazioni di macchina**

$$f(x) = (\tilde{x} \otimes \tilde{x}) \ominus 1 = (\tilde{x} \ominus 1) \otimes (\tilde{x} \oplus 1)$$

Poniamo $g_1(x) = (\tilde{x} \otimes \tilde{x}) \ominus 1$ e $g_2(x) = (\tilde{x} \ominus 1) \otimes (\tilde{x} \oplus 1)$. La formula per l'**errore totale** dell'operazione è

$$\epsilon_{tot1} = \frac{g_1(x) - f(x)}{f(x)}$$

$$\epsilon_{tot2} = \frac{g_2(x) - f(x)}{f(x)}$$

Sviluppiamo $g_1(x)$

$$\begin{aligned} g_1(x) &= ((x(1 + \epsilon_x) \cdot x(1 + \epsilon_x))(1 + \epsilon_1) - 1)(1 + \epsilon_2) \doteq \\ &\doteq (x^2(1 + 2\epsilon_x)(1 + \epsilon_1) - 1)(1 + \epsilon_2) \doteq \\ &\doteq (x^2((1 + 2\epsilon_x + \epsilon_1) - 1)(1 + \epsilon_2) \doteq \\ &\doteq x^2(1 + 2\epsilon_x + \epsilon_1 + \epsilon_2) - (1 + \epsilon_2) = \\ &= (x^2 - 1) + 2x^2\epsilon_x + x^2\epsilon_1 + (x^2 - 1)\epsilon_2 = g_1(x) \\ &|\epsilon_i| \leq u \end{aligned}$$

Quindi, portando al primo membro dell'uguaglianza $\epsilon_{tot1} = \frac{g_1(x) - f(x)}{f(x)}$ tutti i fattori $(x^2 - 1) = f(x)$ si ottiene

$$\epsilon_{tot1} = \frac{2x^2}{x^2 - 1}\epsilon_x + \frac{x^2}{x^2 - 1}\epsilon_1 + \epsilon_2$$

Si evince che l'**errore totale è la somma di due componenti**:

Errore inerente o inevitabile: l'**errore di rappresentazione di x** , che vale 0 se x è numero di macchina ed è **proprietà della funzione**. Nell'esempio, $\epsilon_{in} = \frac{2x^2}{x^2 - 1}\epsilon_x$.
Se l'errore inerente è piccolo si dice che **la funzione è numericamente stabile**, viceversa è **numericamente instabile**.

Errore algoritmico, locale all'operazione e **proprietà dell'algoritmo**. Nell'esempio, $\epsilon_{alg} = \frac{x^2}{x^2 - 1}\epsilon_1 + \epsilon_2$.
Se è piccolo, si dice che l'**espressione è ben condizionata e poco sensibile rispetto alla perturbazione**.
Viceversa, è **mal condizionata**.

Sviluppiamo ora $g_2(x)$ (gli errori δ_i sono analoghi agli ϵ_i , si usa una notazione differente per evidenziare che hanno valori diversi, esseno propri del calcolo e dei valori usati in esso)

$$\begin{aligned} g_2(x) &= ((x(1 + \epsilon_x) - 1)(1 + \delta_1))((x(1 + \epsilon_x) + 1)(1 + \delta_2))(1 + \delta_3) \doteq \\ &\doteq (x^2(1 + \epsilon_x)^2 - 1)(1 + \delta_1 + \delta_2 + \delta_3) \doteq \\ &\doteq (x^2(1 + 2\epsilon_x) - 1)(1 + \delta_1 + \delta_2 + \delta_3) = \\ &= x^2(1 + 2\epsilon_x + \delta_1 + \delta_2 + \delta_3) - (1 + \delta_1 + \delta_2 + \delta_3) = \\ &= (x^2 - 1) + 2x^2\epsilon_x + (x^2 - 1)(\delta_1 + \delta_2 + \delta_3) = g_2(x) \\ &|\epsilon_x| \leq u, |\delta_i| \leq u \end{aligned}$$

Come prima, porto gli $f(x)$ al primo membro dell'uguaglianza $\epsilon_{tot2} = \frac{g_2(x) - f(x)}{f(x)}$ e ottengo

$$\epsilon_{tot2} = \frac{2x^2}{x^2 - 1}\epsilon_x + \delta_1 + \delta_2 + \delta_3$$

Notiamo come

$\epsilon_{in} = \frac{2x^2}{x^2-1}\epsilon_x$, come il calcolo precedente. Infatti, ripetiamo, l'errore inerente è una proprietà della funzione e non di come essa viene calcolata

$\epsilon_{alg} = \delta_1 + \delta_2 + \delta_3$, diverso poiché abbiamo seguito un calcolo differente

Per poter comparare i due algoritmi e scegliere il migliore, analizziamo gli errori algoritmici. L'analisi **va eseguita in valore assoluto**, poiché non si ha alcuna informazione sul segno degli errori, seguendo quindi le regole di comparazione dei valori assoluti:

$$|a + b| \leq |a| + |b|$$

$$|a \cdot b| = |a| \cdot |b|$$

$$\begin{aligned}\epsilon_{alg1} &= \left| \frac{x^2}{x^2-1}\epsilon_1 + \epsilon_2 \right| \leq \left| \frac{x^2}{x^2-1}\epsilon_1 \right| + |\epsilon_2| = \left| \frac{x^2}{x^2-1} \right| \cdot |\epsilon_1| + |\epsilon_2| \leq \frac{x^2}{|x^2-1|}u + u = \left(\frac{x^2}{|x^2-1|} + 1 \right)u \\ \epsilon_{alg2} &= |\delta_1 + \delta_2 + \delta_3| \leq |\delta_1| + |\delta_2| + |\delta_3| \leq 3u\end{aligned}$$

In ϵ_{alg1} , quindi, l'errore algoritmo è minore o uguale a $\left(\frac{x^2}{|x^2-1|} + 1 \right)u$, che però diventa arbitrariamente grande quando x si avvicina ad 1. Il secondo algoritmo è dunque **preferibile**, poiché l'**errore algoritmico è limitato**.

Esempio Calcoliamo ora l'espressione $f(x) = x^2 + 1$. Poniamo $g(x) = (\tilde{x} \otimes \tilde{x}) \oplus 1$ e sviluppiamo, ottenendo

$$g(x) \doteq x^2(1 + 2\epsilon_x + \epsilon_1 + \epsilon_2) + (1 + \epsilon_2)$$

L'errore totale è quindi

$$\epsilon_{tot} = \frac{g(x) - f(x)}{f(x)} = \frac{2x^2}{x^2+1}\epsilon_x + \frac{x^2}{x^2+1}\epsilon_1 + \epsilon_2$$

Studiamo le componenti

$$|\epsilon_{in}| = \left| \frac{2x^2}{x^2+1}\epsilon_x \right| = \left| \frac{2x^2}{x^2+1} \right| \cdot |\epsilon_x| \leq \frac{2x^2}{x^2+1}u \leq 2u$$

perché $\frac{x^2}{x^2+1} \leq 1$. La funzione quindi **non è suscettibile rispetto alle perturbazioni dei dati in ingresso** ed è **ben condizionata**.

$$\epsilon_{alg} = \left| \frac{x^2}{x^2+1}\epsilon_1 + \epsilon_2 \right| \leq \left| \frac{x^2}{x^2+1}\epsilon_1 \right| + |\epsilon_2| = \frac{x^2}{x^2+1}|\epsilon_1| + |\epsilon_2| \leq |\epsilon_1| + |\epsilon_2| \leq 2u$$

Quindi l'algoritmo è **numericamente stabile** perché la componente dell'errore algoritmo è piccola ($\leq 2u$). Questo è il caso migliore che può capitare: **funzione ben condizionata** e **algoritmo numericamente stabile**.

Succede perché non vi è la sottrazione al numeratore. La sottrazione, in macchina, lavora usando molte delle cifre "sporche" della mantissa, cioè quelle che fanno parte del rumore e dell'errore di rappresentazione. Si chiama **errore di cancellazione**.

1.1.3 Errore Inerente

Definizione Si dice **errore inerente** o **inevitabile** generato nel calcolo di $f(x) \neq 0$ la quantità

$$\epsilon_{in} = \frac{f(\tilde{x}) - f(x)}{f(x)}$$

Osserviamo che l'errore inerente è una **proprietà della funzione**, ovvero del problema matematico, e misura la **sensibilità rispetto alla perturbazione del dato iniziale**. Quindi, è indipendente dall'algoritmo di calcolo.

Se l'errore inerente è qualitativamente elevato in valore assoluto, diciamo che **il problema matematico è mal condizionato**, altrimenti si dice che è **ben condizionato**.

1.1.4 Errore Algoritmico

Definizione Si dice **errore algoritmico** generato nel calcolo di $f(x) \neq 0$ la quantità

$$\epsilon_{alg} = \frac{g(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Osserviamo che l'errore algoritmico è **dipendente dall'algoritmo utilizzato**: la funzione g dipende dall'algoritmo usato per calcolare $f(x)$. In generale, **differenti algoritmi conducono a differenti errori algoritmici**.

Se l'errore algoritmico è qualitativamente elevato in valore assoluto si dice che **l'algoritmo è numericamente instabile**, altrimenti si dice che è **numericamente stabile**.

1.1.5 Errore Totale

Definizione Si dice **errore totale** generato nel calcolo di $f(x) \neq 0$ mediante l'algoritmo specificato da g la quantità

$$\epsilon_{tot} = \frac{g(\tilde{x}) - f(x)}{f(x)}$$

L'errore totale misura la **differenza relativa tra l'output atteso e l'output effettivamente calcolato**.

Teorema In un'analisi al primo ordine, vale

$$\epsilon_{tot} = \epsilon_{in} + \epsilon_{alg}$$

Dim

$$\epsilon_{tot} = \frac{g(\tilde{x}) - f(x)}{f(x)} = \frac{f(\tilde{x}) - f(x)}{f(\tilde{x})} \cdot \frac{f(\tilde{x})}{f(x)} + \frac{f(\tilde{x}) - f(x)}{f(x)} = \epsilon_{alg}(1 + \epsilon_{in}) + \epsilon_{in} \doteq \epsilon_{alg} + \epsilon_{in}$$

Viene espresso il fatto che nel calcolo di una funzione razionale, in un'analisi al primo ordine, le due fonti di generazione d'errore forniscono contributi separati che possono essere analizzati indipendentemente.

L'obiettivo dell'analisi numerica è **trovare algoritmi numericamente stabili per problemi ben condizionati**.

1.2 Tecniche per l'Analisi degli Errori

1.2.1 Coefficiente di Amplificazione

Dalla seguente relazione

$$\epsilon_{in} = \frac{f(\tilde{x}) - f(x)}{f(x)} = \frac{f(\tilde{x}) - f(x)}{\tilde{x} - x} \cdot \frac{x}{f(x)} \cdot \frac{\tilde{x} - x}{x}$$

si ricava che la **differenziabilità di $f(x)$ è essenziale** per il controllo dell'errore inerente. Se assumiamo che $f(x)$ è derivabile due volte e continua in (a, b) , allora vale lo sviluppo di Taylor

$$f(\tilde{x}) = f(x) + f'(x)(\tilde{x} - x) + \frac{f''(\xi)}{2}(\tilde{x} - x)^2, \quad |\xi - x| \leq |\tilde{x} - x|$$

da cui si ottiene

$$\epsilon_{in} = \frac{f(\tilde{x}) - f(x)}{f(x)} \doteq \frac{f'(x)}{f(x)} x \epsilon_x = c_x \epsilon_x, \quad c_x = \frac{f'(x)}{f(x)} x$$

La quantità $c_x = c_x(f) = \frac{f'(x)}{f(x)} x$ è detta **coefficiente di amplificazione** e fornisce una **misura del condizionamento del problema**.

In generale, se $f : \Omega \rightarrow R$ è definita su un insieme aperto di R^n , differenziabile due volte su Ω ed il segmento di estremi $\tilde{\mathbf{x}}$ e \mathbf{x} è contenuto in Ω allora vale

$$\epsilon_{in} = \frac{f(\tilde{\mathbf{x}}) - f(\mathbf{x})}{f(\mathbf{x})} \doteq \frac{1}{f(\mathbf{x})} \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\mathbf{x}) x_i \epsilon_{x_i} = \sum_{i=1}^n c_{x_i}(f) \epsilon_{x_i}$$

con

$$c_{x_i}(f) = \frac{1}{f(\mathbf{x})} \frac{\partial f}{\partial x_i}(\mathbf{x}) x_i, \quad 1 \leq i \leq n$$

detti **coefficienti di amplificazione della funzione f rispetto alla variabile x_i** .

Esempio per $f(x) = (x^2 + 1)/x$ si ha

$$c_x = \frac{2 - (x^2 + 1)}{x^2} \cdot \frac{x}{x^2 + 1} \cdot x = \frac{x^2 - 1}{x^2 + 1}$$

e poiché $|c_x| \leq 1$ il problema risulta ben condizionato.

Coefficienti

Per le operazioni aritmetiche si ottiene

$$f(x, y) = x + y \quad \epsilon_{in} \doteq c_x \epsilon_x + c_y \epsilon_y \quad c_x = \frac{x}{x+y} \quad c_y = \frac{y}{x+y}$$

$$f(x, y) = x - y \quad \epsilon_{in} \doteq c_x \epsilon_x + c_y \epsilon_y \quad c_x = \frac{x}{x-y} \quad c_y = -\frac{y}{x-y}$$

$$f(x, y) = x \cdot y \quad \epsilon_{in} \doteq c_x \epsilon_x + c_y \epsilon_y \quad c_x = 1 \quad c_y = 1$$

$$f(x, y) = \frac{x}{y} \quad \epsilon_{in} \doteq c_x \epsilon_x + c_y \epsilon_y \quad c_x = 1 \quad c_y = -1$$

Ne segue, come visto in precedenza, che la sottrazione di due numeri vicini tra loro è potenziale causa di elevata amplificazione degli errori relativi (**cancellazione numerica**). Si nota anche come le operazioni moltiplicative siano ben condizionate.

Esempio Siano $x = 0.2178 \cdot 10^2$ e $y = 0.218 \cdot 10^2$, supponendo di operare con troncamento in base $B = 10$ e $t = 3$ cifre di rappresentazione ($u = 10^{1-t} = 10^{-2}$). Si ha $\tilde{x} = 0.217 \cdot 10^2$ e $\tilde{y} = y$.

Pertanto $\tilde{x} \oplus \tilde{y} = -0.001 \cdot 10^2 = -0.1$ mentre $x - y = -0.0002 \cdot 10^2 = -0.2 \cdot 10^{-1}$ e quindi $|\epsilon_{in}| = \frac{0.8}{0.2} = 0.4$.

Capitolo 2

Problemi dell'Algebra Lineare Numerica

2.1 Norme Matriciali e Vettoriali

I principali problemi dell'algebra lineare concernono la **risoluzione di sistemi lineari** ed il **calcolo di autovalori/autovettori di matrici**. Per studiarne il condizionamento è fondamentale avere degli strumenti per valutare la distanza tra vettori e tra matrici.

La risoluzione di sistemi lineari può essere eseguita in aritmetica reale, ma gli autovalori/autovettori possono essere complessi. Servono quindi **strumenti che operino su spazi vettoriali** F^n e $F^{n \times n}$, con $n \geq 1$ e $F \in \{R, C\}$

Condizionamento su problemi reali $f : R^n \rightarrow R$, $G_m = \frac{f(\tilde{x}) - f(x)}{f(x)}$

Problema: calcolo della soluzione di un sistema lineare

Input $A \in R^{2 \times 2}$, $b \in R^2$

Output $x \in R^2 \mid Ax = b$

Condizionamento: $A, B \rightarrow$ devono essere approssimati a valori di macchina

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \rightarrow \hat{A} = \begin{bmatrix} \hat{a} & \hat{b} \\ \hat{c} & \hat{d} \end{bmatrix} \text{ con } \hat{a} = a(1 + \epsilon_a) \dots$$

$$b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \rightarrow \hat{b} = \begin{bmatrix} \hat{b}_1 \\ \hat{b}_2 \end{bmatrix} \text{ con } \hat{b}_i = b_i(1 + \epsilon_{b_i}) \dots$$

$$Ax = b \rightarrow \hat{A}\hat{x} = \hat{b}$$

Quando \hat{x} è vicino a x ? Per valutare la distanza fra vettori, ho bisogno delle **norme vettoriali**

2.1.1 Norma Vettoriale

Definizione $f : R^n \rightarrow R$ si dice **Norma Vettoriale** se soddisfa le seguenti proprietà

1. $f(x) \geq 0 \wedge (f(x) = 0 \Leftrightarrow x = 0)$
2. $f(\alpha x) = |\alpha|f(x) \quad \forall x \in R^n \quad \forall \alpha \in R$
3. $f(x + y) \leq f(x) + f(y) \quad \forall x, y \in R^n$

Norma Euclidea

$$f(x) = \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

con $x = [x_1, x_2, \dots, x_n]$.

Norma Infinito

$$f(x) = \|x\|_\infty = \max_{i=1}^n |x_i|$$

Dimostrazione che è una norma vettoriale

1. $f(x) = \max |x_i|$ è sempre ≥ 0
 $f(x) = \max |x_i| = 0 \Leftrightarrow |x_i| = 0 \quad \forall i \Leftrightarrow x = 0$
2. $f(\alpha x) = \max |\alpha x_i| = \max |\alpha| |x_i| = |\alpha| \cdot \max |x_i| = |\alpha| f(x)$
3. $f(x + y) = \max |x_i + y_i| \leq \max (|x_i| + |y_i|) \leq \max |x_i| + \max |y_i| = f(x) + f(y)$

Norma Uno

$$f(x) = \|x\|_1 = \sum_{i=1}^n |x_i|$$

In generale le norme di uno stesso vettore hanno valori differenti

Esempio $x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$

$$\|x\|_1 = |2| + |3| = 5$$

$$\|x\|_2 = \sqrt{2^2 + 3^2} = \sqrt{4 + 9} = \sqrt{13}$$

$$\|x\|_\infty = 3$$

Distanza Data $f : R^n \rightarrow R$ **norma**, allora possiamo definire la **distanza su R^n** come

$$dist(x, y) = \|x - y\| = f(x - y)$$

1. $f(x) \geq 0$ e $f(y) \geq 0 \Rightarrow dist \geq 0$ e ($dist = 0 \Leftrightarrow x = y$)
2. $f(\alpha x) = |\alpha| f(x)$
3. $f(x + y) \leq f(x) + f(y)$

Anche la distanza assume valori differenti su norme differenti

Esempio $x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$

$$dist(x, y) = \|x - y\|_1 = \left\| \begin{bmatrix} -1 \\ -3 \end{bmatrix} \right\|_1 = 4$$

$$dist(x, y) = \|x - y\|_2 = \left\| \begin{bmatrix} -1 \\ -3 \end{bmatrix} \right\|_2 = \sqrt{10}$$

$$dist(x, y) = \|x - y\|_\infty = \left\| \begin{bmatrix} -1 \\ -3 \end{bmatrix} \right\|_\infty = 3$$

Equivalenza topologica Può succedere che x e y siano vicini in norma uno e lontani in norma 2? **No** per l'**equivalenza topologica delle norme**: quantitativamente differenti, ma qualitativamente danno gli stessi risultati.

$$\|x - y\|_1 \text{ è piccolo} \Rightarrow \|x - y\|_2 \text{ è piccolo}$$

$$\|x - y\|_\infty \text{ è grande} \Rightarrow \|x - y\|_1 \text{ è grande}$$

2.1.2 Norma Matriciale

Definizione $f : R^{n \times n} \rightarrow R$ si dice **norma matriciale** se soddisfa le seguenti proprietà

1. $f(A) \geq 0 \wedge (f(A) = 0 \Leftrightarrow A = 0 \Leftrightarrow a_{ij} = 0 \ \forall i, j)$
2. $f(\alpha A) = |\alpha| f(A)$
3. $f(A + B) \leq f(A) + f(B)$
4. $f(A \cdot B) \leq f(A) \cdot f(B)$

Difficile definire le norme metriciali, quindi **costruiamo le norme matriciali a partire dalle norme vettoriali**.

Definizione Data $f_V : R^n \rightarrow R$ norma vettoriale, si dice **norma matriciale indotta (compatibile)** la norma $f_M : R^{n \times n} \rightarrow R$ definita da

$$f_M(A) = \max_{\{x \in R^n \mid f_V(x)=1\}} f_V(Ax)$$

1. Prendo i vettori $x \in R^n \mid f_V(x) = 1$, in generale sono infiniti
2. Calcolare Ax , in generale infiniti vettori
3. Prendo $f_V(Ax)$, infiniti **valori**
4. Calcolo il **massimo** di questi valori

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1$$

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty$$

Esempio $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow \|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$.

$\|x\|_2 = 1 \Leftrightarrow \sqrt{x_1^2 + x_2^2} = 1 \Leftrightarrow x_1^2 + x_2^2 = 1$ che è l'equazione per la circonferenza di centro 0 e raggio 1. Quindi tutti i vettori per cui $\|x\|_2 = 1$ sono tutti quelli la cui punta giace sulla circonferenza, quindi tanti quanti i punti della circonferenza **quindi infiniti**.

2.1.3 Teoremi di Caratterizzazione

La definizione quindi **non è pratica** (calcolabile) in generale, perché dovrei fare un massimo su infiniti termini. Quindi si introducono dei **teoremi di caratterizzazione**.

Teorema 1 $\|A\|_2 = \sqrt{\rho(A^T A)}$

Dove $\rho(B) = \max |\lambda_i|$ e $\lambda_i \in \text{spettro}(B)$

$\rho(B)$ è il **raggio spettrale** di B , cioè il **modulo dell'autovalore di modulo massimo** ($\text{spettro}(B)$ è l'insieme degli autovalori).

Teorema 2 $\|A\|_\infty = \max_{i=1}^n \sum_{j=1}^n |a_{ij}|$ **somme per riga poi prendo il massimo**

Teorema 2 $\|A\|_\infty = \max_{j=1}^n \sum_{i=1}^n |a_{ij}|$ **somme per colonna poi prendo il massimo**

Da questi teoremi, se ne deriva che **calcolare la norma uno e la norma infinito è più semplice che calcolare la norma euclidea**.

Vantaggi delle norme matriciali \rightarrow **ulteriorie proprietà** delle norme matriciali indotte

Teorema Sia $\|\bullet\|$ una norma vettoriale e $\|\cdot\|$ la norma matriciale indotta. Allora $\forall A \in R^{n \times n}$, $\forall v \in R^n$ si ha $\|Av\| \leq \|A\| \cdot \|v\|$

Dimostrazione Due casi

$$\begin{aligned} v = 0 \quad \|Av\| &= \|0\| = 0 \\ \|A\| \cdot \|v\| &= \|A\| \cdot 0 = 0 \end{aligned}$$

$$\begin{aligned} v \neq 0 \quad \|Av\| &= \|A \frac{v}{\|v\|} \cdot \|v\|\| \text{ e chiamo il vettore } A \frac{v}{\|v\|} = z. \text{ Quindi, per la seconda propriet\`a avr\`o} \\ &= \|z\| \cdot \|v\| = \|A \frac{v}{\|v\|}\| \cdot \|v\| \leq \|A\| \cdot \|v\| \end{aligned}$$

2.1.4 Condizionamento della risoluzione di un sistema lineare

$$Ax = b \longrightarrow \widehat{A}\widehat{x} = \widehat{b}$$

Per semplicit\`a di analisi, $\widehat{A} = A$ e considero la perturbazione solo su b .

$$Ax = b \longrightarrow A\widehat{x} = \widehat{b}$$

Supponiamo anche A invertibile.

Misurare il condizionamento significa misurare quanto x e \widehat{x} sono vicini, attraverso

$$\epsilon_{in} = \frac{\|x - \widehat{x}\|}{\|x\|}$$

Studiamo la quantit\`a

$$\begin{aligned} \|x - \widehat{x}\| &= \|A^{-1}b - A^{-1}\widehat{b}\| = \\ \widehat{b} &= b + f \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \rightarrow \widehat{b} = \begin{bmatrix} \widehat{b}_1 \\ \widehat{b}_2 \end{bmatrix} = \begin{bmatrix} b_1(1 + \epsilon_1) \\ b_2(1 + \epsilon_2) \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} + \begin{bmatrix} b_1\epsilon_1 \\ b_2\epsilon_2 \end{bmatrix} = b + f \\ &= \|A^{-1}b - A^{-1}(b + f)\| = \|-A^{-1}f\| = \|A^{-1}f\| \leq \|A^{-1}\| \cdot \|f\| \end{aligned}$$

$$\text{Quindi } \|x - \widehat{x}\| \leq \|A^{-1}\| \cdot \|f\|$$

Siccome $\|x\|$ sta al denominatore, maggiorazione a denominatore usa una quantit\`a pi\`u piccola.

$$Ax = b \Rightarrow \|Ax\| = \|b\| \Leftrightarrow \|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$$

$$\text{Quindi } \|x\| \geq \frac{\|b\|}{\|A\|}$$

Da questo ottengo

$$\frac{\|x - \widehat{x}\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|f\|}{\frac{\|b\|}{\|A\|}} = \|A\| \cdot \|A^{-1}\| \cdot \frac{\|f\|}{\|b\|}$$

Quindi

$$\epsilon_{in} = \frac{\|x - \widehat{x}\|}{\|x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|\widehat{b} - b\|}{\|b\|}$$

$\frac{\|\widehat{b} - b\|}{\|b\|}$ **perturbazione relativa sul termine noto**
 $\|A\| \cdot \|A^{-1}\|$ **coefficiente di amplificazione**

Quindi se $\|A\| \cdot \|A^{-1}\|$ \`e grande, allora il sistem \`e mal condizionato.

Se $\|A\| \cdot \|A^{-1}\|$ \`e piccolo, allora \`e ben condizionato. Piccolo quanto?

Osservazione $\|A \cdot A^{-1}\| = \|I\| \leq \|A\| \cdot \|A^{-1}\|$, ma $\|I\| = \max_{\|x\|=1} \|Ix\| = \max_{\|x\|=1} \|x\| = 1$
 Quindi $\|A\| \cdot \|A^{-1}\| \geq 1$ che quindi \`e un coefficiente di amplificazione effettivo, con valore minimo pari a 1.

Capitolo 3

Esercitazioni

3.1 Esercitazione 2

Studiare il **condizionamento** del problema, la **stabilità** degli algoritmi e il **costo computazionale** (stimare il numero di operazioni) di

$$f(x) = x^n = e^{n \log(x)}$$

con $x > 0$

Condizionamento Dato che non è dipendente dal modo in cui lo si calcola, calcoliamo il condizionamento per $f(x) = x^n$ con la formula $\epsilon_{in} = \frac{f'(x)}{f(x)} \cdot x \cdot \epsilon_x$

$$\epsilon_{in} = \frac{n \cdot \cancel{x}^{n-1} \cdot \cancel{x}}{\cancel{x}^n} \epsilon_x = n \epsilon_x$$

quindi

$$|\epsilon_{in}| = |n \cdot \epsilon_x| \leq n \cdot u$$

Dato che $C_x = n$ possiamo dire che il problema è **ben condizionato**

Stabilità Dobbiamo considerare i due algoritmi, perché algoritmi differenti producono errore algoritmico diverso.

$$f(x) = x^n = (((x \cdot x) \cdot x) \cdot x) \dots \cdot x$$

Non possiamo calcolare $x \cdot x$ ma calcoleremo $x \otimes x = x^2(1 + \epsilon_1)$, assumendo $x \in F(B, t, m, M)$ per il calcolo dell'errore algoritmico.

Successivamente calcoliamo $x \cdot x \cdot x$, cioè $(x \otimes x) \otimes x = x^3 \cdot (1 + \epsilon_1) \cdot (1 + \epsilon_2)$.

Quindi $x \cdot \dots \cdot x$ n volte diventerà $x^n \prod_{i=1}^{n-1} (1 + \epsilon_i)$. Quindi

$$g(x) = x^n \prod_{i=1}^{n-1} (1 + \epsilon_i) = x^n (1 + \sum_{i=1}^{n-1} \epsilon_i)$$

$$\epsilon_{alg} = \frac{g(x) - f(x)}{f(x)} = \frac{x^n (1 + \sum_{i=1}^{n-1} \epsilon_i) - x^n}{x^n} = \frac{x^n \sum_{i=1}^{n-1} \epsilon_i}{x^n} = \sum_{i=1}^{n-1} \epsilon_i$$

Quindi

$$|\epsilon_{alg}| = \left| \sum_{i=1}^{n-1} \epsilon_i \right| \leq \sum_{i=1}^{n-1} |\epsilon_i| \leq (n-1)u$$

Quindi l'algoritmo è **numericamente stabile**. In questo caso, ϵ_{in} è dello stesso ordine di ϵ_{alg}

<grafo dell'errore>

L'errore accumulato è $\delta_k = \epsilon_k + \delta_{k-1}$, si somma all'errore locale della moltiplicazione (ϵ_k).

$$|\delta_k| = |\epsilon_k + \delta_{k-1}| \leq |\epsilon_k| + |\delta_{k-1}| \leq |\delta_{k-1}| + u \leq |\delta_{k-2}| + 2u \leq |\delta_1| + (k-1)u \leq |\delta_0| + ku = ku$$

Quindi possiamo dire $|\delta_n| \leq n \cdot u$. Notiamo che gli errori non si amplificano, perché è un'operazione di moltiplicazione, ma vengono semplicemente sommati. La relazione trovata lo dimostra.

Per quanto riguarda il secondo algoritmo

$$f(x) = e^{n \cdot \log(x)}$$

Si osserva che **la funzione è razionale** ($f(x) = x^n$) ma **viene calcolata in macchina usando le funzioni non razionali** $h(x) = e^x$ e $g(x) = \log(x)$. In generale, quindi, non posso eseguire l'analisi dell'errore algoritmico con gli strumenti precedentemente usati. Questo **a meno che non si sappia com'è calcolata la funzione non razionale**. Aggiriamo il problema dicendo che

$$\text{Exp}(x) = e^x \cdot (1 + \epsilon_1)$$

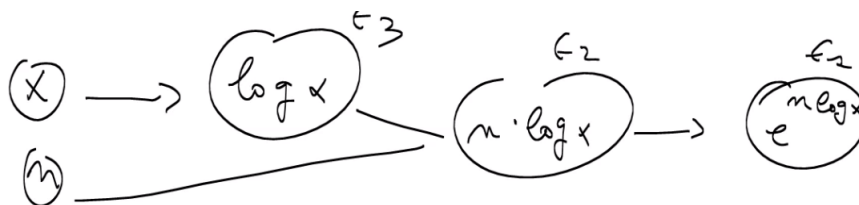
con $|\epsilon_1| \leq u$ e

$$\text{Log}(x) = \log(x) \cdot (1 + \epsilon_2)$$

con $|\epsilon_2| \leq u$. Queste sono approssimazioni, in generale $|\epsilon_1|, |\epsilon_2| \leq k \cdot u$ con k piccola che adesso trascuriamo perché l'analisi qualitativa è analoga.

In questo caso, trattiamo le funzioni non razionali come razionali.

Analisi del grafo



$$|\epsilon_1| \leq u, |\epsilon_2| \leq u, |\epsilon_3| \leq u$$

Per $x \rightarrow \log(x)$ allora $C_x = \frac{\frac{1}{x}}{\log(x)} = \frac{1}{\log(x)}$ ma non amplifica errore su x , perché per il calcolo dell'errore algoritmi assumiamo numeri di macchina.

Per $x \rightarrow e^x$ allora $C_x = \frac{e^x}{e^x} x = x$. Quindi ϵ_2 ha come coefficiente di amplificazione pari a $n \cdot \log(x)$.

Di conseguenza calcolo l'errore algoritmico partendo dalla cima dell'albero e verso le radici (destra verso sinistra)

$$\epsilon_{alg} \doteq \epsilon_1 + (n \cdot \log(x))(\epsilon_2 + \epsilon_3)$$

Concludiamo che $|\epsilon_{alg}| = |\epsilon_1 + n \log(x)(\epsilon_2 + \epsilon_3)| \leq |\epsilon_1| + |n \log(x)(\epsilon_2 + \epsilon_3)| \leq u + |n \log(x)| \cdot |\epsilon_2 + \epsilon_3| \leq u + n |\log(x)| (|\epsilon_2| + |\epsilon_3|) \leq u + 2un |\log(x)| = u \cdot (1 + 2 \cdot n \cdot |\log(x)|)$ cioè

$$|\epsilon_{alg}| \leq u \cdot (1 + 2 \cdot n \cdot |\log(x)|)$$

Concludiamo che l'algoritmo è **numericamente instabile per x piccolo o x è grande**, ma è stabile in un intorno di 1.

L'analisi della stabilità ci dice che il primo algoritmo è preferibile.

Costo computazionale Il primo algoritmo $\Rightarrow (n - 1)$ operazioni moltiplicative, quindi **O(n)**.

Il secondo algoritmo, più complicato perché ci sono due operazioni non razionali oltre alla moltiplicazione.

Esercizio Supposto $n = 2^k$ e $f(x) = x^n = x^{2^k}$, determinare un algoritmo che calcola $f(x)$ con $k = \log_2(n)$ operazioni moltiplicative. Eseguire anche l'analisi dell'errore algoritmico per l'algoritmo trovato.

Capitolo 4

Matlab

11 bit esponente rappresentato in traslazione, 52 bit mantissa (53 cifre rappresentabili per il bit nascosto) $\Rightarrow u = B^{1-t} = 2^{-52}$

4.1 Esponenziale

$f(x) = e^x$ approssimato con Taylor

$$e^x \simeq 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!}$$

Implementiamo un algoritmo che dati x, n in input restituisce $y = 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!}$ in output. Potrebbe essere utile un programma del tipo

```
y = 1;
for k = 1 .. n
    y = y + x^(k)/factorial(k)
end
```

Dal punto di vista del costo computazionale, però, ci sono dei problemi. Al passo k faccio k prodotti per x^k e k prodotti per $k!$, quindi **circa $2k$ prodotti**.

Il costo totale è del tipo $C = 2 + 2 \cdot 2 + 2 \cdot 3 + \dots + 2 \cdot n = 2(1 + 2 + 3 + \dots + n)$.

Ricordiamo $1 + 2 + \dots + n = \frac{n(n+1)}{2}$, quindi $C = 2 \cdot \frac{n(n+1)}{2} = O(n^2)$ quindi costo quadratico.

Una variazione più efficiente

```
y = 1; p = 1; m = 1;
for k = 1 .. n
    p = p * x;      % accumula le potenze
    m = m * k;      % accumula il fattoriale
    y = y + p/m     % accumula la somma
end
```

Al passo k ho 4 operazioni aritmetiche (o 3 operazioni moltiplicative, perché alcune volte si afferma che $t_* \gg t_+$ in termini di operazioni tra bit (**costo booleano**))

Il costo totale $C = 4n$ operazioni (o $3n$).

I due algoritmi possono produrre numeri molto grandi. Una versione che argina questo problema è la seguente.

```
y = 1; z = 1;
for k = 1 .. n
    z = z * x/k;    % controllare l'esplosione dei numeri
    y = y + z;
end
```

Codice MatLab (file: myexp.m)

```
function [y] = myexp(x,n)
    % y = 1 + x + (x^2)/2 + ... + (x^n)/(n!)
    y = 1; z = 1;
    for k = 1:n
        z = z * x/k;
        y = y + z;
    end
end
```