

Name: Sam

MATH-UA 252-005 / MA-UY.4424-C - Midterm #1

1. Newton's method and minimization:

1. Show how to use Newton's method to minimize a function  $f(x)$ .
2. Prove that this is equivalent to minimizing a sequence of quadratic polynomials. (Hint: use a Taylor expansion to find the quadratic polynomial to be minimized.)

1. Let  $x^*$  be minimizing argument of  $f$ .

we require  $f'(x^*) = 0$ . ("first-order necessary conditions for optimality")

2. Let  $f(x_n + \Delta x) = f(x_n) + f'(x_n)\Delta x + \frac{1}{2}f''(x_n)\Delta x^2 + O(\Delta x^3)$   
define  $p_n(\Delta x) = f(x_n) + f'(x_n)\Delta x + \frac{1}{2}f''(x_n)\Delta x^2$ .

$$p'_n(\Delta x) = f'(x_n) + f''(x_n)\Delta x = 0$$

$$\Rightarrow \Delta x = -\frac{f'(x_n)}{f''(x_n)}$$

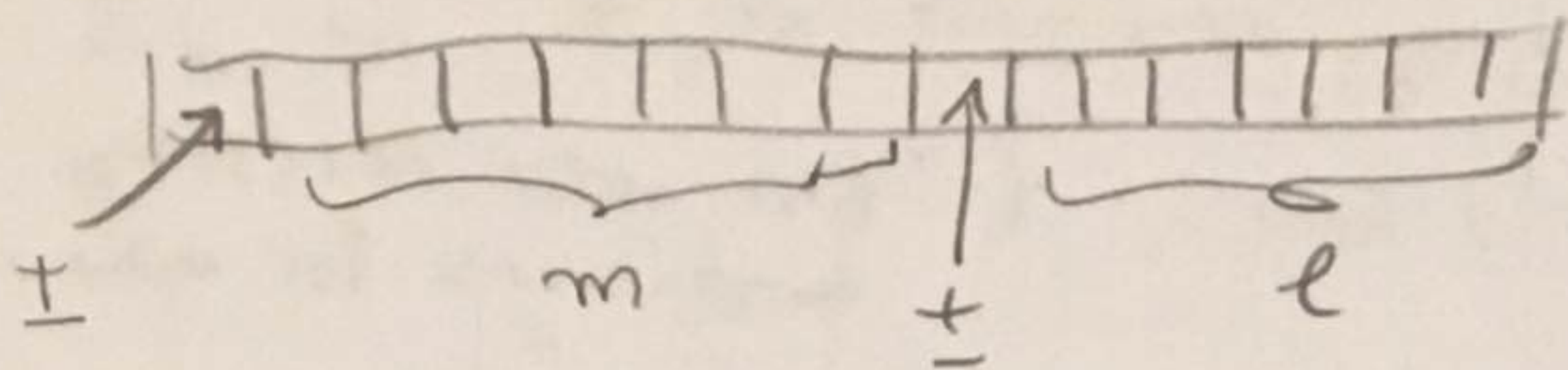
If we set  $x_{n+1} = x_n + \Delta x = x_n - \frac{f'(x_n)}{f''(x_n)}$ ,

we just get the Newton iteration for solving  $f'(x) = 0$ .

2. Floating-point numbers:

1. Make up and explain a simple floating-point representation which uses 16 bits.
2. Let  $x \in [0, 1]$ . We can write  $x = t_{-1}3^{-1} + t_{-2}3^{-2} + t_{-3}3^{-3} + \dots$  for some coefficients  $t_k$ , where  $t_k \in \{0, 1, 2\}$ . We call each coefficient  $t_k$  a *trit*. Find this expansion for  $x = 1/10$ . (Hint: each  $t_k$  must be nonnegative!)

1.  $\pm m \times 2^{\pm e}$



2.  $\frac{1}{10} = \frac{1}{1+3^2} = \frac{1}{3^2} \cdot \frac{1}{1+3^{-2}}$

geometric series  $\rightarrow = \frac{1}{3^2} \left( 3^0 - 3^{-2} + 3^{-4} - 3^{-6} + 3^{-8} - 3^{-10} + \dots \right)$

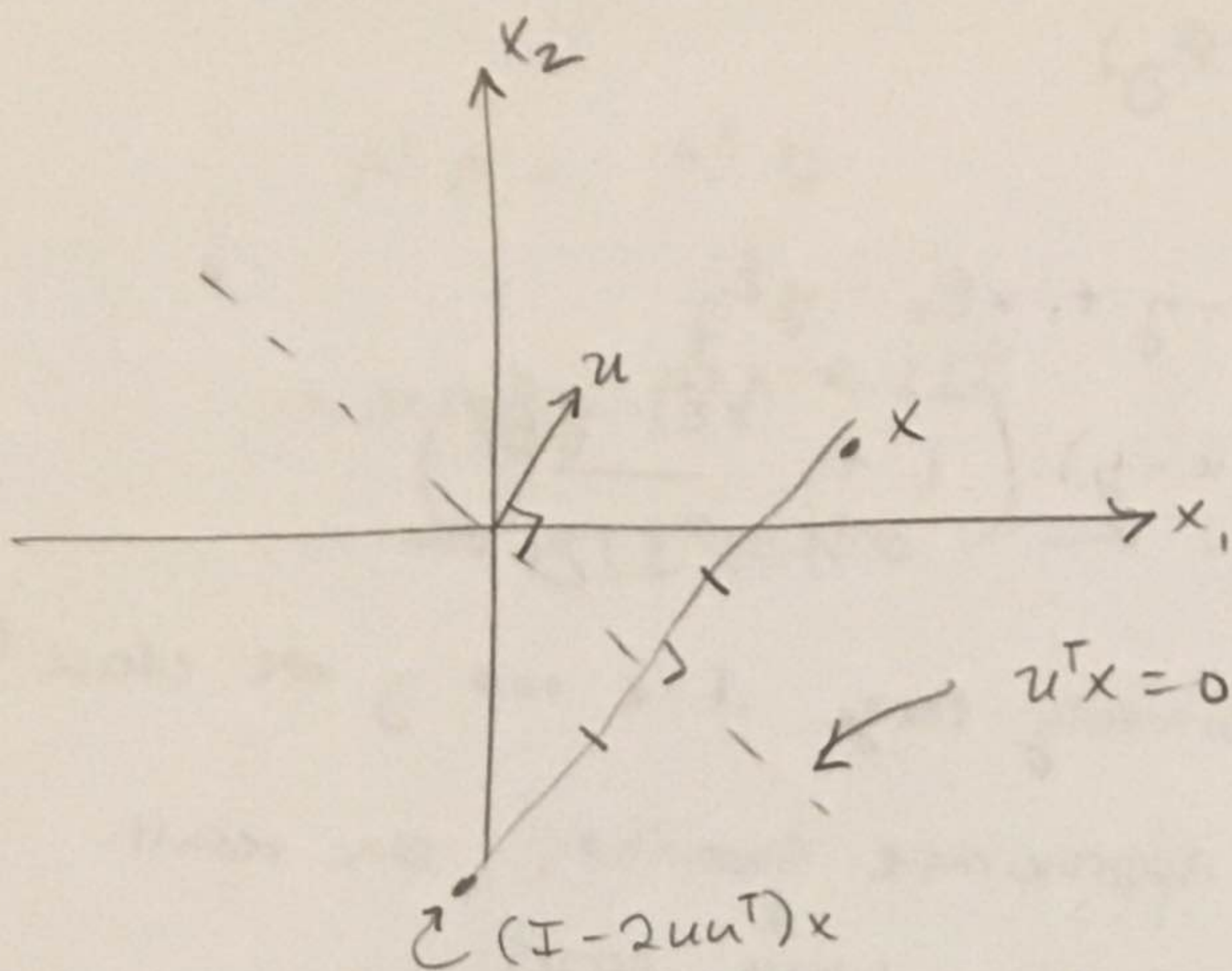
$\downarrow 3^0 - 3^{-2} = 1 - \frac{1}{9} = \frac{10}{9} = 3^{-2} (2 \cdot 3^1 + 2 \cdot 3^0) = 2 \cdot 3^{-1} + 2 \cdot 3^{-2}$

$= 3^{-2} \cdot \left( 2 \cdot 3^{-1} + 2 \cdot 3^{-2} + 2 \cdot 3^{-5} + 2 \cdot 3^{-6} + \dots \right)$

$= 0.00\overline{22}00\overline{22}00\overline{22}00\dots_3$

### 3. Linear algebra.

- Let  $u \in \mathbb{R}^n$  be a unit vector, and let  $R_n = I - 2uu^T$ . Draw a picture for  $n = 2$  which shows what multiplying with  $R_n$  does. For general  $n$ , compute the eigenvalues of  $R_n$  (including their multiplicity), explain what the eigenspaces are (the subspaces of  $\mathbb{R}^n$  corresponding to each distinct eigenvalue), and compute  $\det(R_n)$ .
- Prove that if  $A$  is positive definite, then  $a_{ii} > 0$  for each  $i$ .



$$2. \quad a_{ii} = e_i^T A e_i > 0.$$

Eigenvalues: two cases:

$$\begin{aligned} \text{i.) } u^T x = 0 &\Rightarrow (I - 2uu^T)x = x \\ &\Rightarrow \text{eigenvalue} = 1 \end{aligned}$$

$$\begin{aligned} \text{ii.) } u^T x \neq 0 & \Rightarrow u^T x = u \\ &\Rightarrow (I - 2uu^T)u = u - 2|u|^2 u \\ &= (1 - 2|u|^2)u \\ &\Rightarrow \text{eigenvalue} = 1 - 2|u|^2 \end{aligned}$$

The first eigenvalue has multiplicity  $n-1$ , and the second has multiplicity  $1$ , since the dimension of the plane  $u^T x = 0$  is  $n-1$ .

This describes the two eigenspaces.

Now we can conclude that  $\det(I - 2uu^T) = 1 - 2|u|^2$ .

4. Sources of error. Compute the relative error of subtraction and explain the phenomenon of *catastrophic cancellation*.

$$\hat{x} = x(1 + \epsilon_x)$$

$$\hat{y} = y(1 + \epsilon_y)$$

$$\begin{aligned}\hat{x} - \hat{y} &= x - y + x\epsilon_x - y\epsilon_y \\ &= (x - y) \left( 1 + \frac{x\epsilon_x - y\epsilon_y}{x - y} \right)\end{aligned}$$

This can be arbitrarily large if  $x$  and  $y$  are close!

So subtracting approximate quantities can result in arbitrarily large relative errors.

This is "catastrophic cancellation."

5. Numerical linear algebra. Let  $Ax = b$  be an overdetermined linear system. Give a quick explanation of how to use the Cholesky decomposition to solve this overdetermined system. What must be true of  $A$  for this to work? Why?

Normal equations:

$$A^T A x = A^T b$$

cholesky:  $A^T A = LL^T$

$$\Rightarrow LL^T x = A^T b \rightarrow L^T x = \underbrace{L^{-T} L^{-1} A^T b}_{\text{mult } (O(mn))}$$

$\underbrace{\hspace{10em}}_{\text{tri solve } (O(n^2))}$   
 $\underbrace{\hspace{10em}}_{\text{tri solve } (O(n^2))}$

Cholesky decomposition exists for spd matrices.

1)  $A^T A$  is symmetric ( $\underline{s}$  pd) ✓

2) if  $A$  has full column rank, then  $A^T A$  is positive definite ( $\underline{s}$  pd) ✓

Why? IF  $A$  is not full column rank, must be a vector  $u$  such that  $Au = 0$ . Hence  $A^T A u = 0$ . So  $A^T A$  has at least one zero eigenvalue. So  $A^T A$  cannot be pd.

But observe that  $z^T A^T A z = \|Az\|^2 \geq 0$  for any  $z \neq 0$ . So  $A^T A$  is at least positive semi def. Hence, full column rank  $\Rightarrow A^T A$  pos def.

6. **The Babylonian algorithm.** The iteration for the Babylonian algorithm is  $x_{n+1} \leftarrow (y/x_n + x_n)/2$ , where  $y = x^2$ . We saw that if  $y > 0$  and  $x_0 > 0$ , then  $x_n \rightarrow \sqrt{y} = x$ . Consider the error  $e_n = x_n - x$ . The sequence  $x_n \rightarrow x$  with an order of convergence equal to  $q > 0$  if:

$$\lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n^q} = \mu \in [0, 1]. \quad (1)$$

If we pick  $q = 1$  and find that the limit above equals  $\mu = 0$ , then the sequence converges *superlinearly*. Assume that  $x_n \rightarrow \sqrt{y} = x$  and show that the Babylonian algorithm converges superlinearly.

Note:

$$\frac{1}{2} \left( \frac{x^2}{x_n} + x_n \right) - x = \frac{1}{2} \left( \frac{x^2}{x_n} + \frac{x_n^2}{x_n} - \frac{2xx_n}{x_n} \right)$$

$$= \frac{1}{2} \left( \frac{(x - x_n)^2}{x_n} \right)$$

and

So:

$$\frac{e_{n+1}}{e_n} = \frac{1}{2} \frac{(x_n - x)^2}{x_n} \cdot \frac{1}{x_n - x} = \frac{1}{2} \frac{x_n - x}{x_n}$$

$$= \frac{1}{2} \left( 1 - \frac{x}{x_n} \right)$$

Hence:

$$\lim_{n \rightarrow \infty} \frac{e_{n+1}}{e_n} = \lim_{n \rightarrow \infty} \frac{1}{2} \left( 1 - \frac{x}{x_n} \right)$$

$$= \frac{1}{2} \left( 1 - \frac{x}{x} \right) = 0,$$

since  $x_n \rightarrow x$  as  $n \rightarrow \infty$ !