

SECTION C

NUMERICAL METHODS 1 (3.09) taught by David Ham

Candidates being examined in “Numerical Methods 1” should answer at least one question from Section C.

- C1. (i) Consider the number 5.25
(a) Write the number in the form:

$$s1.m \times 2^{e-b} \quad (1)$$

Use a floating point format with 3 exponent bits and 4 mantissa bits, and a bias of 3. All the numbers must be written in binary.
(6 marks)

- (b) Convert this number into the bit pattern which would actually be stored, according to the layout provided on the formula sheet.
(2 marks)

- (ii) (a) The Bisection method is given by the algorithm:

```
fl ← f(xl)
repeat
  xc ←  $\frac{x_l + x_r}{2}$ 
  fc ← f(xc)
  if fcfl > 0 then
    xl ← xc
    fl ← fc
  else
    xr ← xc
until |f(xc)| < ε
```

Produce a sketch of one iteration of this algorithm for the function $f(x) = x^2 - 1$ with a starting interval $x_l = 0.5, x_r = 2$. Show all the relevant points and lines and indicate which end of the interval will be removed.
(3 marks)

- (b) The Newton-Raphson iteration is given by the algorithm:

```
repeat
  fx ← f(x)
  x ←  $x - \frac{f_x}{f'(x)}$ 
until (|fx| < ε) or maximum iterations exceeded.
```

Produce a sketch of one iteration of this algorithm for the function $f(x) = x^2 - 1$ with a starting position $x_0 = 0.5$. Show all the relevant points and lines.
(3 marks)

- (c) What is the key advantage of Newton-Raphson iteration over the bisection method?
(1 mark)
- (d) What additional information is required to use Newton-Raphson iteration rather than the bisection method?
(1 mark)

(e) What advantage does the bisection method have over Newton-Raphson iteration?

(1 mark)

(iii) Suppose that A is a rectangular $n \times m$ matrix with $n > m$, and that the columns of A are all orthogonal to each other, and each column has a modulus of 1 ($|A_{:,i}| = \sqrt{A_{:,i} \cdot A_{:,i}} = 1$ for $0 \leq i \leq m - 1$).

(a) prove that $A^T A = I$ where I is the $m \times m$ identity matrix.

(4 marks)

(b) prove therefore that the least squares solution to the overdetermined matrix equation:

$$A\mathbf{x} = \mathbf{b}$$

is given by $\mathbf{x} = A^T \mathbf{b}$

(4 marks)

C2. (i) Convert the numbers in the following problems into 4 bit two's complement signed binary integers before performing the calculation and converting back into base 10.

(a) $-5 + 3$

(4 marks)

(b) -1×4

(4 marks)

(ii) For each of the following Python functions, write a mathematical expression using only matrices and vectors which performs the same operation. In each case, state whether each input and output is a matrix or vector.

```
(a) def function_1(a,b,c):
    import numpy

    d=numpy.zeros(numpy.size(b))

    for j in range(len(d)):
        d[j]=numpy.dot(a[j,:],c)-b[j]

    return d
```

(4 marks)

```
(b) def function_2(a):
    import numpy

    c=a.shape[1]
    b=numpy.zeros((c,c))

    for i in range(c):
        for j in range(c):
            b[i,j]=numpy.dot(a[:,i],a[:,j])

    return b
```

(4 marks)

(iii) A central difference approximation to the third derivative of a function $f(x)$ is given by:

$$f'''(x) = \frac{-f(x-2h) + 2f(x-h) - 2f(x+h) + f(x+2h)}{2h^3} + \mathcal{O}(h^p)$$

Where h is some small positive step size, and p is the order of convergence. Using Taylor series for the function f , or otherwise, prove that the order of convergence, p , is equal to 2.

You may use without proof any of the properties of \mathcal{O} .

(9 marks)