537 Exam 1

Fall 2025

Tuesday, September 30, 2025

Name:	Solutions
-------	-----------

Person number:

(a) (b) (c) (c) (d) (d) (d) (d) (e) (e) (e) (e) (e) (e) (e) (e) (e) (e	①①①②③④(6)(7)(8)(9)	①①①②③④(5)(6)(7)(8)(9)	①①①②③④(5)(6)(7)(8)(9)	(a)(b)(c)(d)(d)(e)(e)(f)(e)(f)(f)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)(g)<l< td=""><td>①①①②③④(5)(6)(7)(8)(9)</td><td>① (1)② (3)④ (4)⑤ (6)(7)(8)(9)</td><td>(a) (b) (c) (</td></l<>	①①①②③④(5)(6)(7)(8)(9)	① (1)② (3)④ (4)⑤ (6)(7)(8)(9)	(a) (b) (c) (



537 Fall 2025

EXAM INSTRUCTIONS

The exam is closed-book. No calculator.

One single side of handwritten notes is permitted.

- Relax.
- ANSWER ANY 4 OF THE 5 QUESTIONS
- Work steadily and carefully.
- You have 20 minutes per question.
 If you get bogged down on a question, move on.
- For maximum credit, show all your work.
- Do not waste time writing information that is not asked for.
- About 5% will be awarded for "style": be sure your <u>writing is easy to read</u>, and make your reasoning, explanations, and calculations <u>clear</u>, <u>explicit</u>, <u>easy to follow</u>.

1

(a) i. How many distinct floating-point numbers are there in the IEEE 754 64-bit system?ii. What fraction of the numbers do we lose by reserving the top exponent for non-numbers (nan, inf)?

Briefly explain your answers.

iii. Also put the following numbers in increasing order, using the fact that $2^p \approx 10^{0.30 p}$:

S = the number of stars in a typical galaxy $\sim 10^{11}$

M = the number of molecules in a drop of water $\sim 5 \times 10^{20}$

N = the number of distinct floating-point numbers in the IEEE 754 64-bit system.

(b) What is the **definition** and the **significance** of the number ϵ_{mach} for a floating-point arithmetic system.

Emach is the spacing between I and the next machine number. It measures the resolution of the system of machine numbers. It is the maximum relative spacing between (non-zero) machine numbers, and twice the maximum relative error in rounding a real number (in range) to a machine number.

IMPORTANT: there is a part (c) on the next page.

(c) Small integers are represented exactly in IEEE 754 64-bit floating-point format.

(Recall that this format has a sign bit, an 11-bit exponent biased by $1023 = 2^{10} - 1$, and a 52-bit mantissa.) Discover what is smallest positive integer that <u>cannot</u> be represented exactly. Explain your reasoning clearly,

and also give the decimal order of magnitude of the answer.

For economy suppose for a moment we have a 4-bit mantissa (instead of \$2-bit).

Then the small positive integers are represented as follows:

1 1.0000 × 2° = 10

2 1.0000 × 2' = 10

3 1.1000 × 2' = 11

$$2^{4+1}$$
 1.1111 × 2⁴ = 11111

 2^{4+1} 1.0000 × 2⁵ = 100000

All these integers so far are machine numbers.

However,

 2^{4+1} = 100001

would require an additional bit in the mantissa to represent exactly.

Thus for a 52-bit mantissa, the answer is 2^{52+1} = 2^{53} + 1 \approx $10^{15.9}$ \approx 10^{16} .

Check:

\$./hexcodeout

Enter number: 9007199254740991

433ffffffffffffff

Enter number: 9007199254740992

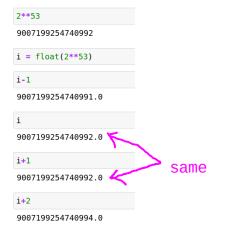
43400000000000000

Enter number: 9007199254740993

4340000000000000

Enter number: 9007199254740994

4340000000000001



The pink text is just for motivation: you can completely ignore it for now. In a round of "musical chairs", let *p* denote the number of players, and *c* the number of available chairs. If the players choose seats at random, the expected number of players surviving to the next round is

$$\left(1-e^{-\frac{p}{c}}\right)c$$

(a) Explain why direct evaluation of this formula in finite-precision arithmetic is problematic for 0 .

This can be relevant if c, the number of "chairs", is very large, as might be the case if the "players" are insect parasites, and "chairs" are the hosts they parasitize.

For $O , <math>e^{-pc} \approx 1$, so we are subtracting near-equals which we know results in severe loss of significance.

- (b) How would you evaluate this function accurately when $p \le c$? Be specific. If you introduce any truncation error, it should be $O(p^3)$ or smaller.
- Use Taylor approx for $e^{-P/c}$: $e^{\times} = 1 + \times + \frac{x^2}{2!} + O(x^3)$ $e^{-P/c} = 1 \frac{2}{c} + \frac{p^2}{2c^2} + O((\frac{p_c}{c})^3)$ and $(1 e^{-P/c})c = c(\frac{p^2}{c} \frac{p^2}{2c^2} + O((\frac{p_c}{c})^3))$ $= p \frac{p^2}{2c} + O((\frac{p_c}{c})^3)$ (Contextual sanity check: when there are many more "chairs" than players, almost everyone survives: #survivors $\approx p$.)

Find the approximate maximum relative error in the floating point evaluation of this expression (x+c)-c

when x is small and not necessarily a machine number, and c is a machine number greater than 1.

Answer in terms of |x|, c, and $\varepsilon_{\text{mach}}$, and just keep the term or terms that dominate as |x| goes to zero.

Suggestion: to avoid mistakes, expand everything fully before canceling anything.

| teletive error | = |
$$\frac{f(x)}{(x(1+S_1)+C)(1+S_2)} - C \frac{1}{(1+S_3)} - x$$
 | | $\frac{1}{(x+S_2)} + xS_1S_2 + cS_2 - \frac{1}{(x+S_3)} + xS_1S_2 + cS_2 - \frac{1}{(x+S_3)} + xS_1S_2 + cS_2 + xS_1S_2 + cS_2S_3 + xS_1S_2S_3 + xS_1S_2S_3 + cS_2S_3 + cS_2S_3$

4

Dividing by multiplying

Suppose you have a computer chip that can do multiplication but not division. Consider working around the deficit by obtaining the reciprocal of a real number $r \neq 0$ by finding a root of an appropriate function using Newton's method.

Each of the following two functions has a root at 1/r:

(i)
$$f_1(x) = rx - 1$$

(ii)
$$f_2(x) = 1 - \frac{1}{rx}$$

- (a) Which would be satisfactory for the purpose described above? Explain in detail.
- (b) For the one that works, what is the order of convergence? Justify your answer by citing an appropriate theorem.

(a) Newton iteration is $X_{K+1} = X_K - \frac{f(X_K)}{f'(X_K)}$ (i) $f_i(x) = r$ So $\times_{K+1} = \times_K - (\Gamma \times_{K-1}) = \frac{1}{\Gamma}$ So Newton converges to immediately. But division is required to compute f. So not useful. (ii) f2(x)=+1/2, SO XK+1 = XK - (1-1/XK) = XK-LXX+XK = 2xx-rx2. This can be computed without division! (b) f₂(t)= r ≠ 0 and f₂ ∈ C'(Rt). So by Thm 2.4 or Thm 2.5 Newton on fz converges (at least) gnadratically.

Test Newton on f 2:

```
def reciprocal(r):
    # crude starting guesses
    if r>1:
        x = .5
    else:
        x = 2
    tol = le-15 # relative error tolerance
    i = 0
    print(i,x)
    while True:
        s = x - r*x*x
        x += s
        i += 1
        print(i,x)
        if abs(r*s) <= tol: return x
        if i>8:
             print('convergence not achieved')
             break
```

```
reciprocal(3)

0 0.5
1 0.25
2 0.3125
3 0.33203125
4 0.3333282470703125
5 0.333333333335557231
6 0.3333333333333333
7 0.3333333333333333
```

```
reciprocal(.1)

0 2
1 3.6
2 5.904
3 8.3222784
4 9.718525023289343
5 9.992077183748574
6 9.999993722898264
7 9.99999999999606
8 9.999999999999999
9 10.0
```

5

Convergence of functional iteration

Suppose a function $g: \mathbb{R} \to \mathbb{R}$ has the following properties:

- g(z) = z,
- g is continuously differentiable on an open interval containing z, and
- |g'(z)| < 1.

Use the Mean Value Theorem to prove that there is an interval I about z such that g is a contraction on I and g maps I into itself (so that the Contraction Mapping Theorem guarantees convergence to z for a sufficiently close starting point).

Since
$$|g'(z)| < 1$$
 and g' is critical g' is critical g' is critical g' is continuous, we can choose an $\varepsilon > 0$ (small enough) such that on the interval $I=[z-\varepsilon,z+\varepsilon]$, $|g'| \le L < 1$.

Then for $\chi \text{Im } I$, the Mean Value Theorem says $|g(x)-g(y)| < |g'(\xi)(x-y)| \le L|x-y|$.

That is, g' is L -Lipschitz with $L < 1$, which is the definition of a critical g' .

Therefore g' is f in f in