

# 2025 年度 科学技術計算

学生番号 36714143  
情報工学科知能情報分野  
和田小百合

2025 年 11 月 4 日

## 目次

1	課題 02-2	2
1.1	課題の内容	2
1.2	IEEE 754 binary64 とは	2
1.3	コードの説明	2
1.4	a + b の計算	5
2	課題 2-05	6
2.1	課題の内容	6
2.2	softmax 関数とは	6
2.3	愚直に実装した softmax 関数の実装方法と結果	6
2.4	演習資料で説明されていた数値的に安定な softmax 関数の実装方法と結果	7
3	課題 2-06	8
3.1	課題の内容	8
3.2	実装方法について	8
3.3	実行結果	8

## 1 課題 02-2

### 1.1 課題の内容

コンピュータで計算をするときは 2 進数に変換する際に誤差を生んでしまうことがあります。これについて bit 表現から解くことによってなぜ誤差が含まれてしまうか説明します。

### 1.2 IEEE 754 binary64 とは

IEEE 754 binary64 とは、2 進数浮動小数点形式の 1 種です

構成としては、

- 一般形： $(-1)^{\text{sign}}(d_0.d_1d_2d_3 \cdots d_{52})_2 \times 2^{n - 1023}$ 
  - sign：符号 1 ビット（正負を表す）
  - $(d_0.d_1d_2d_3 \cdots d_{52})_2$ ：仮数部 52 ビット + 正規化表現の仮定された 1 ビットで計 53 ビットの 2 進数
    - \* 正規化表現の場合、 $d_0 = 1$  が仮定される（格納されない）
    - \*  $n = (n_{11} \cdots n_2 n_1 n_0)_2$ ：指数部 11 ビット ( $1 \leq n \leq 2046$ , バイアスは 1023)
    - \*  $n$  は符号なし整数で表現され、指数部の範囲は  $-1022 \leq n - 1023 \leq 1023$  となる。

### 1.3 コードの説明

```

1 a = 18.0
2 b = 2.0**57
3 c = (a + b) - b
4 print(c)

```

Listing 1 課題 02-2

上のコードは本来  $c = 18.0$  になるはずですが、実際には 32.0 が出力されます。この理由を浮動小数点数の演算手順をビット単位で追跡することでなぜ違う値が出力されたを説明します。まず、 $a = 18.0$ ,  $b = 2.0^{57}$  を IEEE754 binary64 のビット表現に変換し、それぞれの符号・指数・仮数を確認する。

また、 $a$ ,  $b$  の符号・指数・仮数部分を表示するために今回の講義の教科書部分にあるコードを用い、表示させました。

```

1 def float64_to_binary(f: float | np.float64, separate: bool = False) -> str | Tuple[str, str, str]:
2     """convert float64 to 64 bits as string
3
4     Args:
5         f (float | np.float64): IEEE754 binary64 floating point number
6         separate (bool, optional): return tuple of (sign bit, exponent bits,
7             mantissa bits) if True,
8             return 64 bits if False. Defaults to False.
9
10    Returns:
11        str or Tuple of str: bit string of f
12    """
13    assert isinstance(f, float)
14    from struct import pack, unpack

```

```

14
15     # see https://note.nkmk.me/python-float-hex/
16     s = format(unpack('>Q', pack('>d', f))[0], "064b")
17
18     if separate:
19         return s[0:1], s[1:12], s[12:]
20     else:
21         return s
22
23 def sign_str_to_str(s: str) -> str:
24     """sign bit to +/-
25
26     Args:
27         s (str): "0" or "1"
28
29     Returns:
30         str: "+" if s is "0" else "-" (if s is "1")
31     """
32     assert isinstance(s, str)
33     assert len(s) == 1
34     assert s in ["0", "1"]
35
36     return "+" if s == "0" else "-"
37
38 def exponent_str_to_int(
39     e: str,
40     len_e: int = 11,
41     bias: int = 1023
42 ) -> Tuple[int, bool]:
43     """exponent string to exponent integer
44
45     Args:
46         e (str): exponent string consisting of "0" and "1"
47         len_e (int, optional): the number of bits for the exponent. Defaults to
48             11.
49         bias (int, optional): bias to exponent. Defaults to 1023.
50
51     Returns:
52         Tuple[int, bool]: exponent in decimal integer, and True if normalized
53             number else False
54     """
55     assert isinstance(e, str)
56     assert len(e) == len_e
57     for bit in e:
58         assert bit in ["0", "1"]
59
60     p = sum([2**i * int(bit) for i, bit in enumerate(e[::-1])])
61     if p > 0:
62         p -= bias
63         is_normalized = True
64     else:
65         p = -bias + 1
66         is_normalized = False

```

```

65     if p > bias:
66         p = np.nan
67     return p, is_normalized
68
69
70 def mantissa_str_to_float(
71     m: str,
72     is_normalized: bool,
73     len_m: int = 52
74 ) -> float:
75     """mantissa string to mantissa integer
76
77     Args:
78         m (str): mantissa string consisting of "0" and "1"
79         is_normalized (bool): True if m is normalized number else False
80         len_m (int, optional): the number of bits for the mantissa. Defaults to
81             52.
82
83     Returns:
84         float: mantissa in decimal float
85     """
86     assert isinstance(m, str)
87     assert len(m) == len_m
88     for bit in m:
89         assert bit in ["0", "1"]
90
91     p = 1 if is_normalized else 0
92     p += sum([2 ** -(i + 1) * int(bit) for i, bit in enumerate(m)])
93
94     return p
95
96 def sem_to_str(s: str, e: str, m: str, n_spaces: int = 0) -> str:
97     """convert output strings of 'float64_to_binary()' with 'separate=True' to
98     a human-readable string
99
100    Args:
101        s, e, m (str): bit strings of sign, exponent, and mantissa.
102        n_spaces (int, optional): the number of leading white spaces on the
103            left. Defaults to 0.
104
105    Returns:
106        str: human-readable string representation of a number represented by s,
107            e, and m
108
109    """
110    s = sign_str_to_str(s)
111    e, is_normalized = exponent_str_to_int(e)
112    m = mantissa_str_to_float(m, is_normalized)
113
114    n_spaces = " " * n_spaces
115
116    if e is not np.nan:
117        return f"{n_spaces}{s} {2^{e:+7d}} * {m}"
118    else:

```

```
return f"{n_spaces}{s} nan/inf    * {m}"
```

Listing 2 課題 02-2/a b を bit で表示する関数

そして、この関数を用い、それぞれの文字が IEEE754 binary64 でどのように表示されているか確認するコードを以下に示します。

```
1 print(f"decimal    s exponent      mantissa")
2 for f in [
3     a,
4     b,
5     a + b,
6     -b,
7     c,
8 ]:
9     s, e, m = float64_to_binary(f, separate=True)
10    print(f"{f:.2f} --> {s} {e} {m}")
11    print(sem_to_str(s, e, m, n_spaces=9))
```

そして、実行結果が以下になります。

## 1.4 a + b の計算

$a + b$  の計算手順は以下です。

1. 指数の整列
  2. 假数の加算
  3. 正規化
  4. 丸め

この処理手順に従って説明する。

1. 指数の整列まず、 $a, b$  は表によって示されているとおりです。 $a$  の桁数は指数部分より、 $2^4$  で、 $b$  は  $2^{57}$  である。仮数を計算するためにまず桁を同じにする必要があり、桁が大きい方すなわち  $b$  のほうに合わせる。桁数の合わせ方は、仮数を右に何 bit 右シフトすることで桁を合わせる。今回の場合 52bit 右シフトする必要があります。右シフトをした  $a$  (表では  $a'$  とする) はこのようになります。

2. 仮数の加算桁数が同じになったため仮数を足し合わせる。また、丸みを考慮するため 52bit 以降の桁も記します。

- 正規化正規化は、一番左の整数部分が 1 ではないときに行うものである。今回は整数部分が 1 であるため特に何も行わない。
  - 丸め丸めでは飛び出てしまった部分が 0.5 以上であればくりあげを行う。今回の場合は繰り上げられるため、結果は  $(a+b)'$  で表示する。この時に、 $b$  が非常に大きいため、 $a$  の寄与が落ちてしまうのである。

## 2 課題 2-05

## 2.1 課題の内容

softmax 関数をコンピュータで実装するときに、コンピュータの計算の限界から、実際の値との誤差が生まれることを検証しました。

## 2.2 softmax 関数とは

softmax 関数とは、シグモイド関数を多次元に拡張した関数であり、多クラス分類問題において、ニューラルネットワークの出力を確率分布(0,1)に変換することができる。機械学習や深層学習で使用される関数である。softmax 関数  $\sigma(x)$  は次のように定義される。

$$\sigma(x)_i = \frac{e^{x_i}}{\sum_{j=1}^M e^{x_j}}$$

### 2.3 愚直に実装した softmax 関数の実装方法と結果

### 2.3.1 softmax 関数の実装方法

上式のまま python で実装したものが以下のコードである。

```
1 #愚直に実装した数値的に不安定なsoftmax
2 def unstable_softmax(x: np.ndarray,) -> np.ndarray:
3     a = np.sum(np.exp(x))
4     y = np.exp(x) / a
5
6     return y
```

Listing 3 課題 02-05/愚直に実装した softmax 関数

コードを見てみると、softmax 関数には  $\exp$  をデータ数に応じて足し合わせる動作があるため、簡単にオーバーフローしてしまうことが考えられます。

### 2.3.2 結果

結果の出力には以下のコードを使用した。

```

1  rng = np.random.default_rng()
2  a = rng.integers(800, 1000, size=10)
3  a = a.astype(np.float64)
4  print(a.dtype)
5
6  y = unstable_softmax(a)
7  print(y)
8  for i in y:
9      print(f"{i:.15e}")

```

NumPy に含まれる乱数生成の default\_rng でジェネレータのインスタンスを作り、インスタンスに対して、確率分布メソッドを呼び出し、型を決定して softmax 関数に入れました。実行結果は nan になり、オーバーフローしていることがわかりました。

## 2.4 演習資料で説明されていた数値的に安定な softmax 関数の実装方法と結果

### 2.4.1 演習資料で説明された softmax 関数の実装方法

愚直に実装した softmax 関数ではオーバーフローしてしまったため、softmax 関数並進不変性を利用する。

$$\begin{aligned}\sigma(x+y)_i &= \frac{e^{x_i+y}}{\sum_{j=1}^M e^{x_j+y}} \\ &= \frac{e^y e^{x_i}}{e^y \sum_{j=1}^M e^{x_j}} \\ &= \sigma(x)_i\end{aligned}$$

また、最大値  $\tilde{x} = \max_i x_i$  を用いて、log-sum-exp trick を適用します。この式を実装したものが以下のコードになります。

```

1 def stable_softmax(x: np.ndarray) -> np.ndarray:
2     max = np.max(x)
3     exp_x = np.exp(x - max)
4     sum_exp_x = np.sum(exp_x)
5
6     y = np.exp(x - max - np.log(sum_exp_x)) # log
7
8
9     return y

```

Listing 4 課題 02-05/演習資料で説明されていた数値的に安定な softmax 関数

### 2.4.2 結果

以下が結果の一例である。

```

1.587178613846118e-60
4.157641866918001e-07
4.999997921178430e-01
3.833822442893757e-53
1.546674362601999e-50
1.652849626184181e-37
1.271832294905762e-13

```

```
1.172775181673195e-59
4.999997921178430e-01
5.056105360903163e-44
```

この結果から softmax 関数がオーバーフローしないことがわかりました。また、scipy.special.softmax の結果が以下のようになります。

```
array([1.58717861e-60, 4.15764187e-07, 4.99999792e-01, 3.83382244e-53, 1.54667436e-50, 1.65284963e-37,
1.27183229e-13, 1.17277518e-59, 4.99999792e-01, 5.05610536e-44])
```

より、結果が一致するため数値が安定的であるといえる。

### 3 課題 2-06

#### 3.1 課題の内容

$1 + \epsilon > 1$  となる（つまり演算の前後で値が変化する）ような最小の浮動小数点数  $\epsilon$  マシンイプシロンの値を求めるアルゴリズムを作成します。

#### 3.2 実装方法について

以下のように関数にまとめ、どの型でも反映できるようにしました。

```
1 def machine_epsilon(dtype):
2     eps = dtype(1)
3     one = dtype(1)
4
5     while one + eps > one:
6         eps = eps / dtype(2)    # εを半分にする
7         print(f"{eps:.15e}")    # 表記で表示
8
9     eps = eps * dtype(2)    # 最後に倍して戻す
10    print("Machine epsilon =", f"{eps:.15e}")
11
12    return eps
```

基本的なアルゴリズムのまま作成しているため、詳細なコードの説明は省略する。

#### 3.3 実行結果

```
アルゴリズムで計算したマシンイプシロンと np.finfo で出力した悔過を以下にしめす。np.float64
Machine epsilon = 2.220446049250313e-16
np.float32
Machine epsilon = 1.192092895507812e-07
np.float16
Machine epsilon = 9.765625000000000e-04
```

```
np.finfo()
2.220446049250313e-16
1.1920929e-07
0.000977
```

このように結果が一致することがわかりました。